

CS 78 Computer Networks

Transport Layer

Andrew T. Campbell

<http://www.cs.dartmouth.edu/~cs78/>

These slides come with the course book Computer Networks and are reproduced under the authors' copyright
All material copyright 1996-2006, J.F. Kurose and K.W. Ross, All Rights Reserved

Transport Layer 3-1

Chapter 3: Transport Layer

Our goals:

- r understand principles behind transport layer services:
 - m multiplexing/demultiplexing
 - m reliable data transfer
 - m flow control
 - m congestion control
- r learn about transport layer protocols in the Internet:
 - m UDP: connectionless transport
 - m TCP: connection-oriented transport
 - m TCP congestion control

Transport Layer 3-2

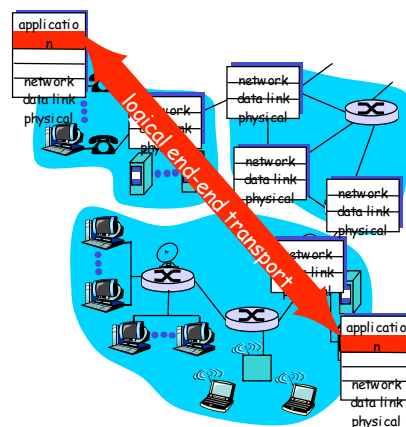
Chapter 3 outline

- r 3.1 Transport-layer services
- r 3.2 Multiplexing and demultiplexing
- r 3.3 Connectionless transport: UDP
- r 3.4 Principles of reliable data transfer
- r 3.5 Connection-oriented transport: TCP
 - m segment structure
 - m reliable data transfer
 - m flow control
 - m connection management
- r 3.6 Principles of congestion control
- r 3.7 TCP congestion control

Transport Layer 3-3

Transport services and protocols

- r provide *logical communication* between app processes running on different hosts
- r transport protocols run in end systems
 - m send side: breaks app messages into *segments*, passes to network layer
 - m rcv side: reassembles segments into messages, passes to app layer
- r more than one transport protocol available to apps
 - m Internet: TCP and UDP



Transport Layer 3-4

Transport vs. network layer

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

Household analogy:

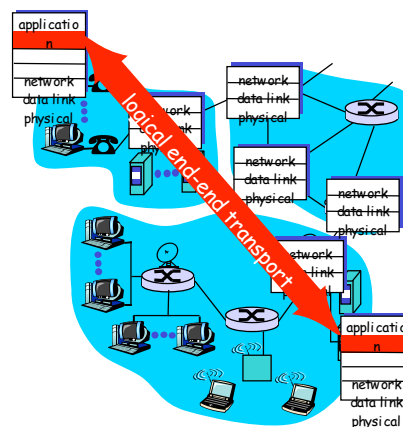
12 kids sending letters to 12 kids

- r processes = kids
- r app messages = letters in envelopes
- r hosts = houses
- r transport protocol = Ann and Bill
- r network-layer protocol = postal service

Transport Layer 3-5

Internet transport-layer protocols

- r reliable, in-order delivery (TCP)
 - m congestion control
 - m flow control
 - m connection setup
- r unreliable, unordered delivery: UDP
 - m no-frills extension of "best-effort" IP
- r services not available:
 - m delay guarantees
 - m bandwidth guarantees



Transport Layer 3-6

Chapter 3 outline

- r 3.1 Transport-layer services
- r 3.2 Multiplexing and demultiplexing
- r 3.3 Connectionless transport: UDP
- r 3.4 Principles of reliable data transfer
- r 3.5 Connection-oriented transport: TCP
 - m segment structure
 - m reliable data transfer
 - m flow control
 - m connection management
- r 3.6 Principles of congestion control
- r 3.7 TCP congestion control

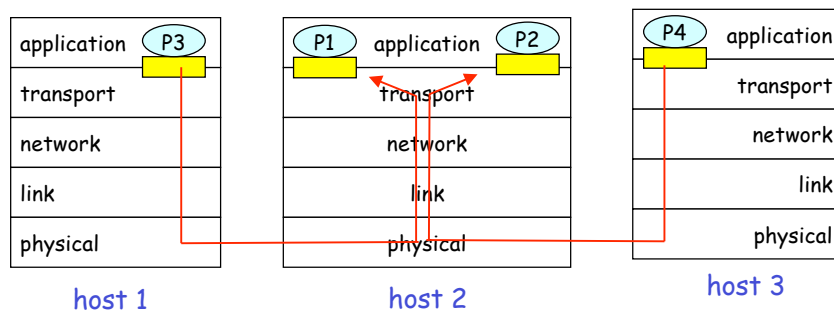
Transport Layer 3-7

Multiplexing/demultiplexing

Demultiplexing at rcv host:
delivering received segments to correct socket

Multiplexing at send host:
gathering data from multiple sockets, enveloping data with header (later used for demultiplexing)

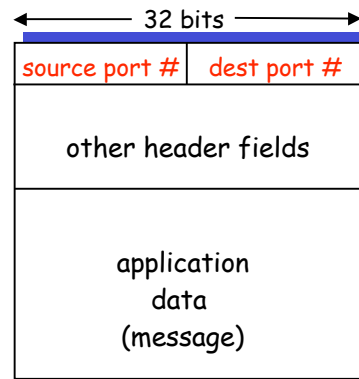
■ = socket ○ = process



Transport Layer 3-8

How demultiplexing works

- r host receives IP datagrams
 - m each datagram has source IP address, destination IP address
 - m each datagram carries 1 transport-layer segment
 - m each segment has source, destination port number
- r host uses IP addresses & port numbers to direct segment to appropriate socket



TCP/UDP segment format

Transport Layer 3-9

Connectionless demultiplexing

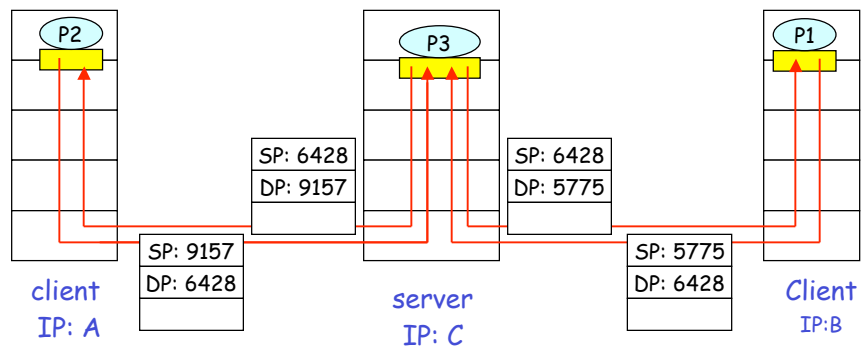
- r Create sockets with port numbers:

```
DatagramSocket mySocket1 = new
    DatagramSocket(12534);
DatagramSocket mySocket2 = new
    DatagramSocket(12535);
```
- r UDP socket identified by two-tuple:
(dest IP address, dest port number)
- r When host receives UDP segment:
 - m checks destination port number in segment
 - m directs UDP segment to socket with that port number
- r IP datagrams with different source IP addresses and/or source port numbers directed to same socket

Transport Layer 3-10

Connectionless demux (cont)

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



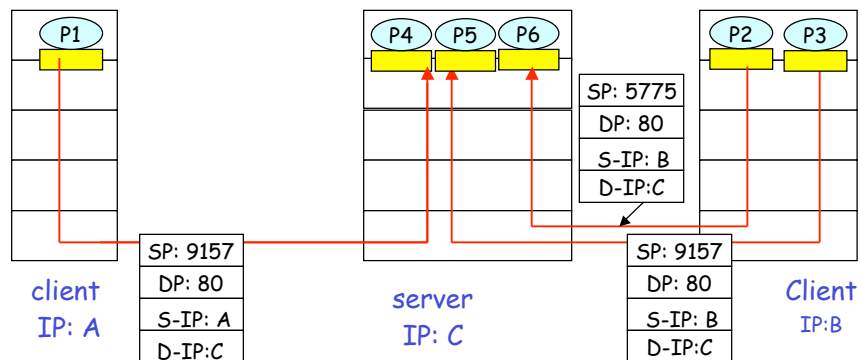
Transport Layer 3-11

Connection-oriented demux

- r TCP socket identified by 4-tuple:
 - m source IP address
 - m source port number
 - m dest IP address
 - m dest port number
- r recv host uses all four values to direct segment to appropriate socket
- r Server host may support many simultaneous TCP sockets:
 - m each socket identified by its own 4-tuple
- r Web servers have different sockets for each connecting client
 - m non-persistent HTTP will have different socket for each request

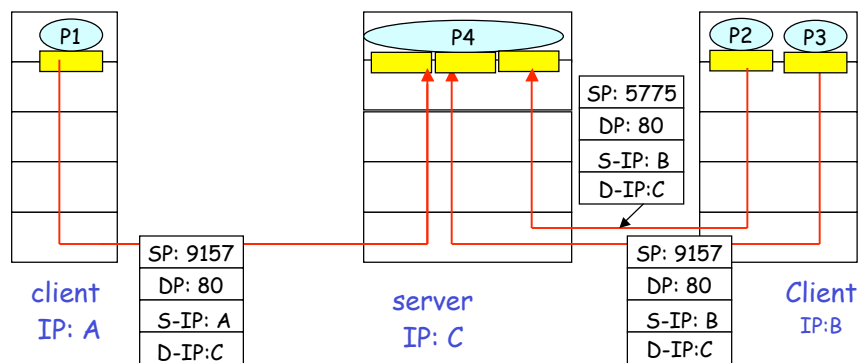
Transport Layer 3-12

Connection-oriented demux (cont)



Transport Layer 3-13

Connection-oriented demux: Threaded Web Server



Transport Layer 3-14

Chapter 3 outline

- r 3.1 Transport-layer services
- r 3.2 Multiplexing and demultiplexing
- r 3.3 Connectionless transport: UDP
- r 3.4 Principles of reliable data transfer
- r 3.5 Connection-oriented transport: TCP
 - m segment structure
 - m reliable data transfer
 - m flow control
 - m connection management
- r 3.6 Principles of congestion control
- r 3.7 TCP congestion control

Transport Layer 3-15

UDP: User Datagram Protocol [RFC 768]

- r "no frills," "bare bones" Internet transport protocol
- r "best effort" service, UDP segments may be:
 - m lost
 - m delivered out of order to app
- r **connectionless:**
 - m no handshaking between UDP sender, receiver
 - m each UDP segment handled independently of others

Why is there a UDP?

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

Transport Layer 3-16

UDP: more

- r often used for streaming multimedia apps

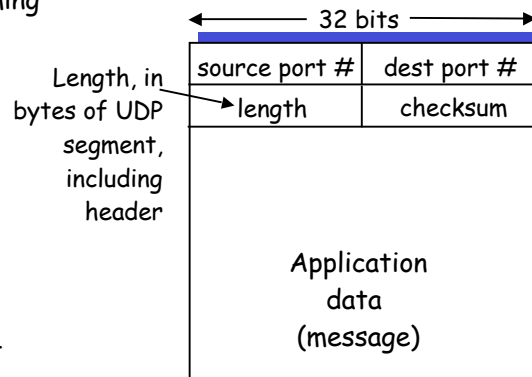
- m loss tolerant
- m rate sensitive

- r other UDP uses

- m DNS
- m SNMP

- r reliable transfer over UDP: add reliability at application layer

- m application-specific error recovery!



UDP segment format

Transport Layer 3-17

UDP checksum

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

Sender:

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

Receiver:

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

....

Transport Layer 3-18

Internet Checksum Example

- r Note

- m When adding numbers, a carryout from the most significant bit needs to be added to the result

- r 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
sum																
checksum	1	0	1	1	1	0	1	1	1	0	1	1	1	1	0	0
	0	1	0	0	0	1	0	0	0	1	0	0	0	0	1	1

Transport Layer 3-19

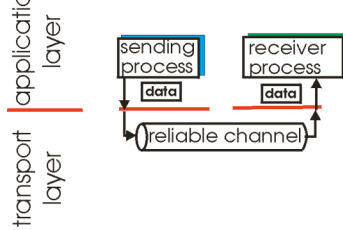
Chapter 3 outline

- r 3.1 Transport-layer services
- r 3.2 Multiplexing and demultiplexing
- r 3.3 Connectionless transport: UDP
- r 3.4 Principles of reliable data transfer
- r 3.5 Connection-oriented transport: TCP
 - m segment structure
 - m reliable data transfer
 - m flow control
 - m connection management
- r 3.6 Principles of congestion control
- r 3.7 TCP congestion control

Transport Layer 3-20

Principles of Reliable data transfer

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



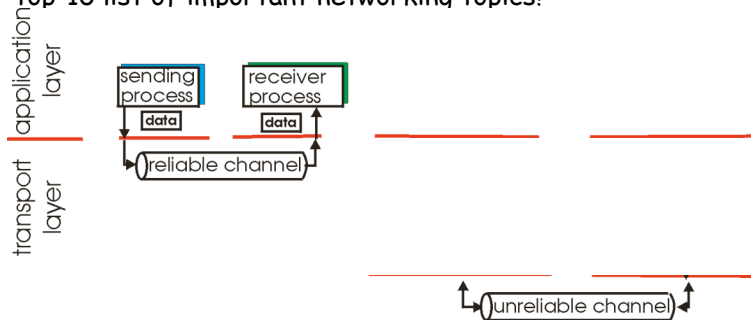
(a) provided service

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

Transport Layer 3-21

Principles of Reliable data transfer

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



(a) provided service

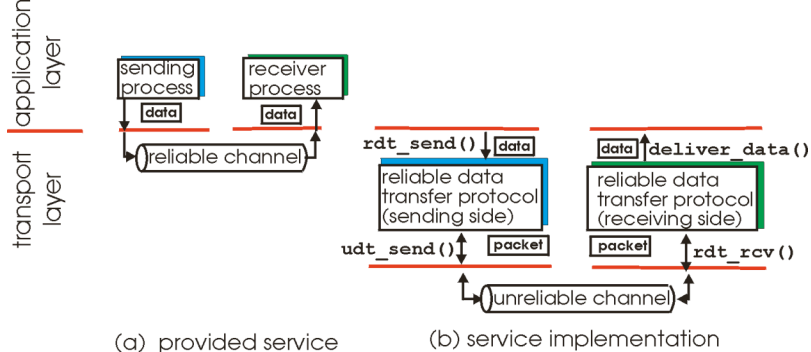
(b) service implementation

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

Transport Layer 3-22

Principles of Reliable data transfer

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



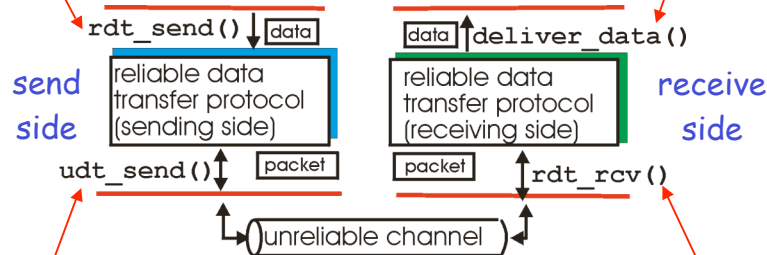
- r characteristics of unreliable channel will determine complexity of reliable data transfer protocol (rdt)

Transport Layer 3-23

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



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

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

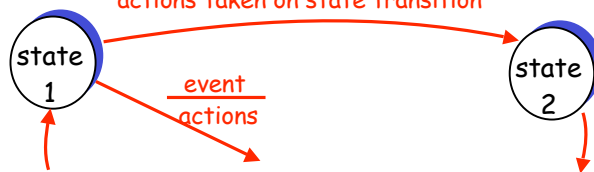
Transport Layer 3-24

Reliable data transfer: getting started

We'll:

- r incrementally develop sender, receiver sides of reliable data transfer protocol (rdt)
- r consider only unidirectional data transfer
 - m but control info will flow on both directions!
- r use finite state machines (FSM) to specify sender, receiver

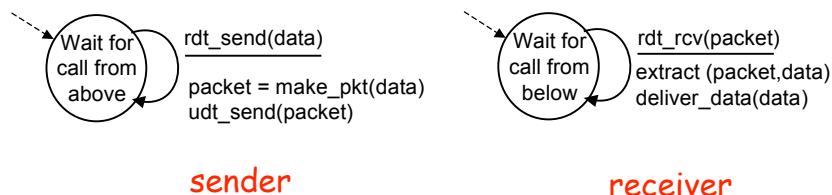
state: when in this "state" next state uniquely determined by next event



Transport Layer 3-25

Rdt1.0: reliable transfer over a reliable channel

- r underlying channel perfectly reliable
 - m no bit errors
 - m no loss of packets
- r separate FSMs for sender, receiver:
 - m sender sends data into underlying channel
 - m receiver read data from underlying channel



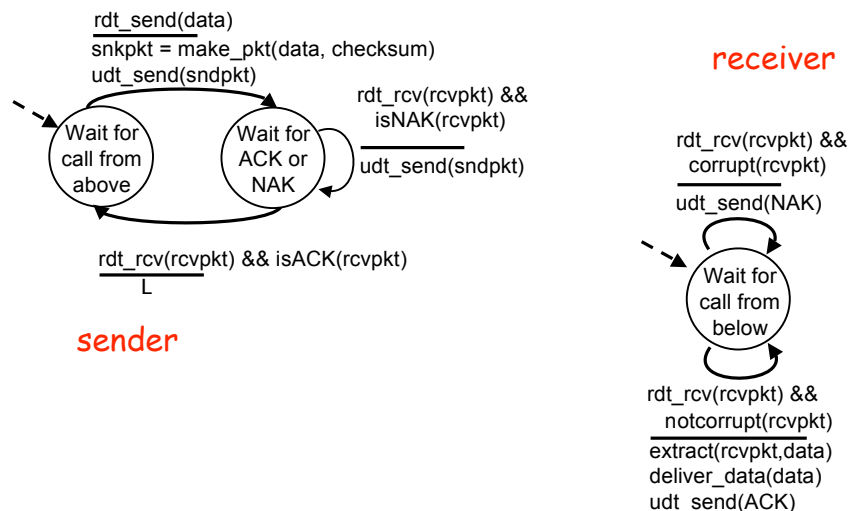
Transport Layer 3-26

Rdt2.0: channel with bit errors

- r underlying channel may flip bits in packet
 - m checksum to detect bit errors
- r *the question: how to recover from errors:*
 - m *acknowledgements (ACKs)*: receiver explicitly tells sender that pkt received OK
 - m *negative acknowledgements (NAKs)*: receiver explicitly tells sender that pkt had errors
 - m sender retransmits pkt on receipt of NAK
- r new mechanisms in rdt2.0 (beyond rdt1.0):
 - m error detection
 - m receiver feedback: control msgs (ACK,NAK) rcvr->sender

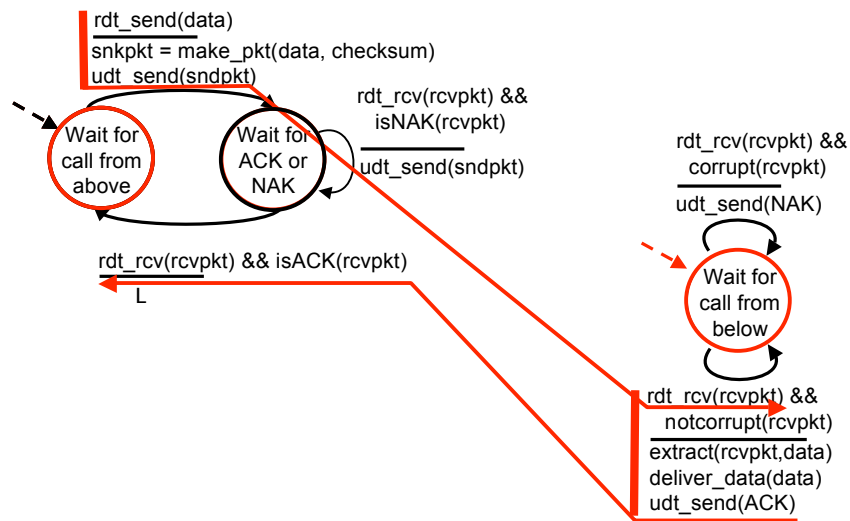
Transport Layer 3-27

rdt2.0: FSM specification



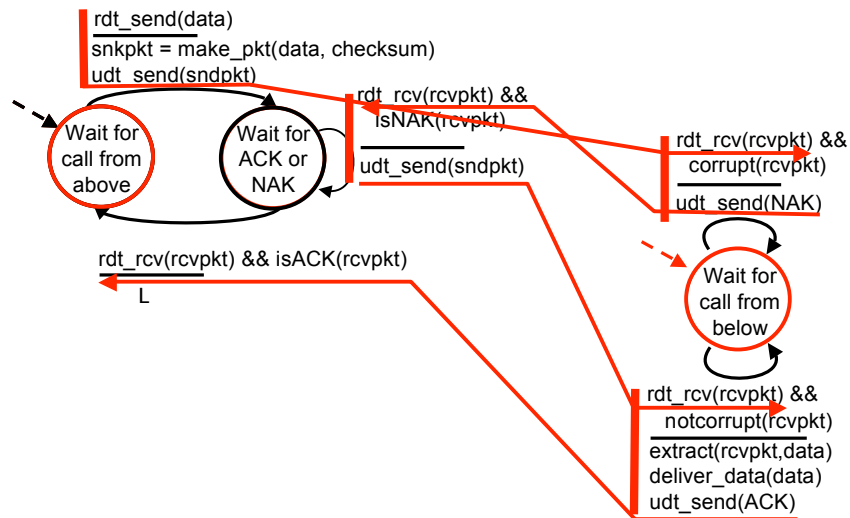
Transport Layer 3-28

rdt2.0: operation with no errors



Transport Layer 3-29

rdt2.0: error scenario



Transport Layer 3-30

rdt2.0 has a fatal flaw!

What happens if ACK/NAK corrupted?

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

stop and wait

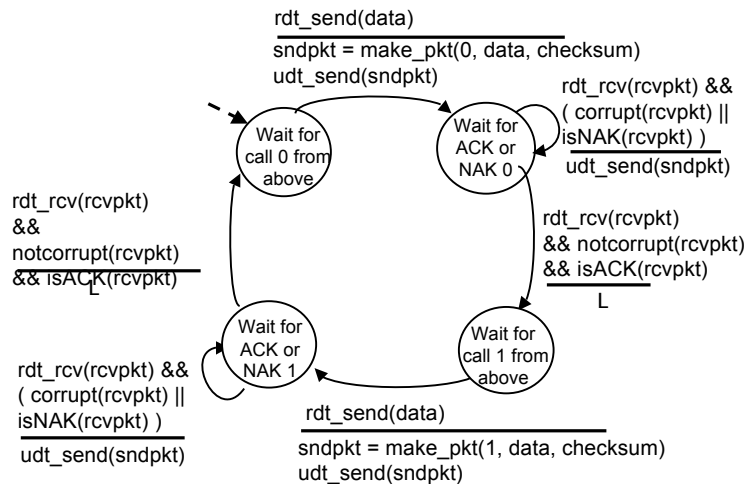
Handling duplicates:

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

Sender sends one packet, then waits for receiver response

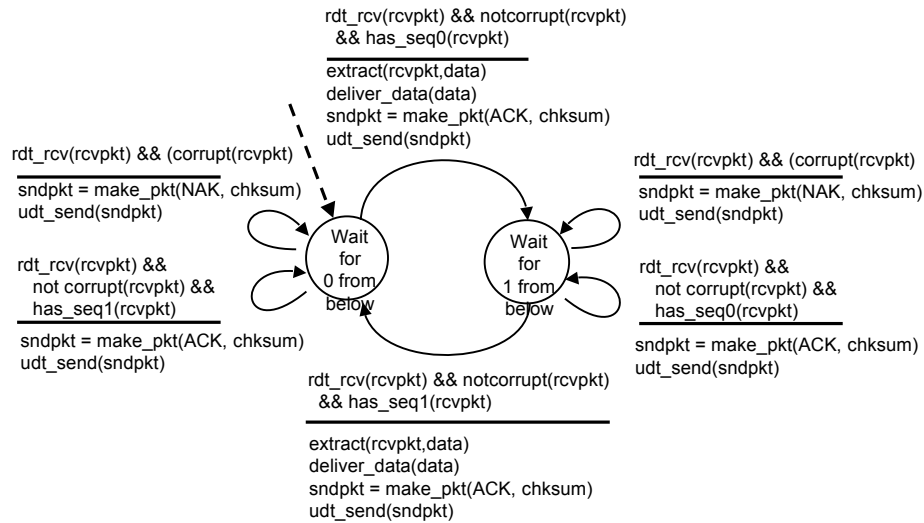
Transport Layer 3-31

rdt2.1: sender, handles garbled ACK/NAKs



Transport Layer 3-32

rdt2.1: receiver, handles garbled ACK/NAKs



Transport Layer 3-33

rdt2.1: discussion

Sender:

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

Receiver:

- r must check if received packet is duplicate
 - m state indicates whether 0 or 1 is expected pkt seq #
- r note: receiver can *not* know if its last ACK/NAK received OK at sender

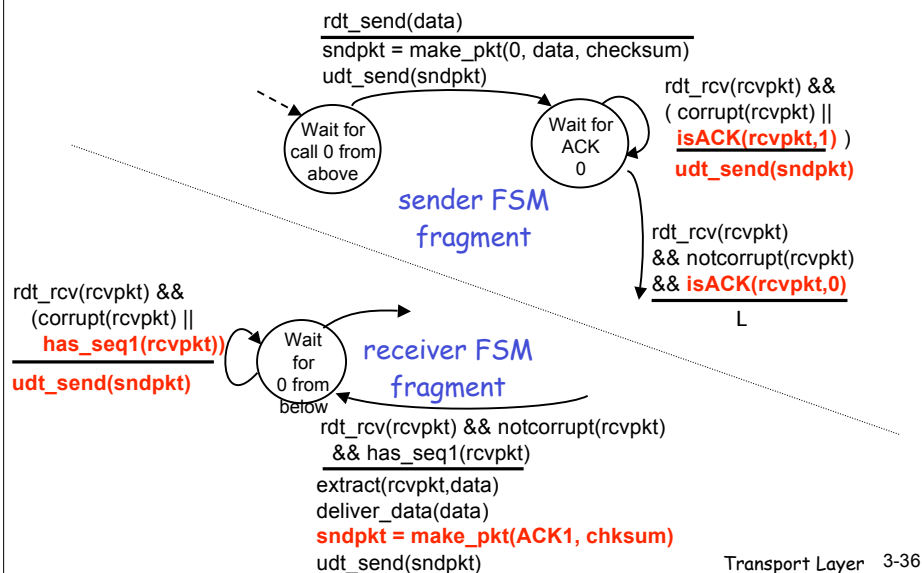
Transport Layer 3-34

rdt2.2: a NAK-free protocol

- r same functionality as rdt2.1, using ACKs only
- r instead of NAK, receiver sends ACK for last pkt received OK
 - m receiver must *explicitly* include seq # of pkt being ACKed
- r duplicate ACK at sender results in same action as NAK: *retransmit current pkt*

Transport Layer 3-35

rdt2.2: sender, receiver fragments



rdt3.0: channels with errors and loss

New assumption:

underlying channel can
also lose packets
(data or ACKs)

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

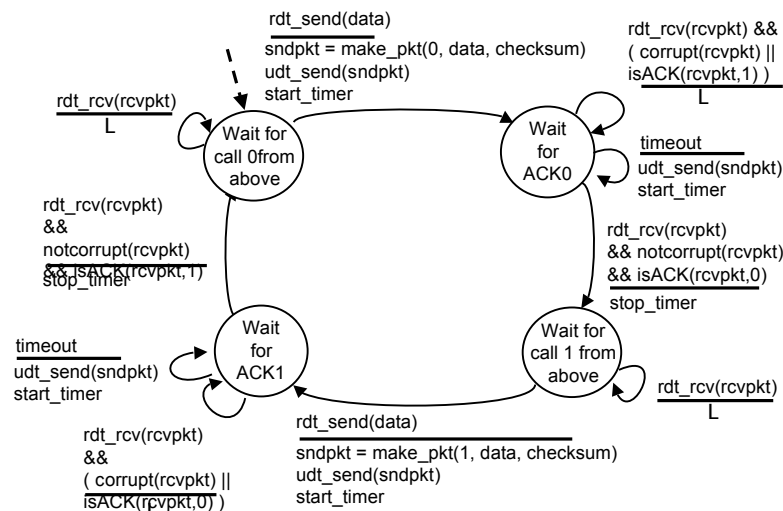
Approach: sender waits

"reasonable" amount of time
for ACK

- r retransmits if no ACK received
in this time
- r if pkt (or ACK) just delayed (not
lost):
 - m retransmission will be
duplicate, but use of seq. #'s
already handles this
 - m receiver must specify seq #
of pkt being ACKed
- r requires countdown timer

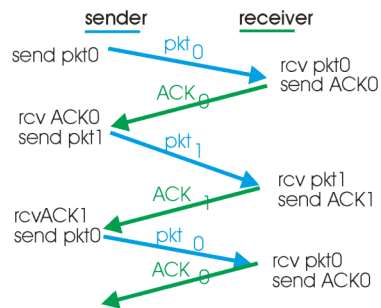
Transport Layer 3-37

rdt3.0 sender

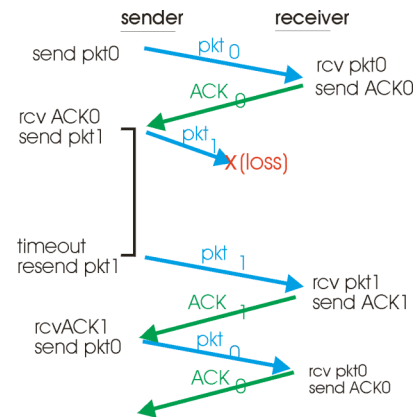


Transport Layer 3-38

rdt3.0 in action



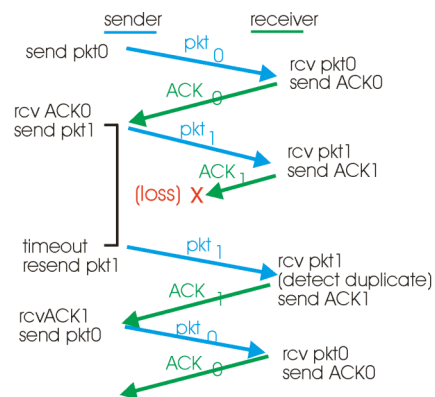
(a) operation with no loss



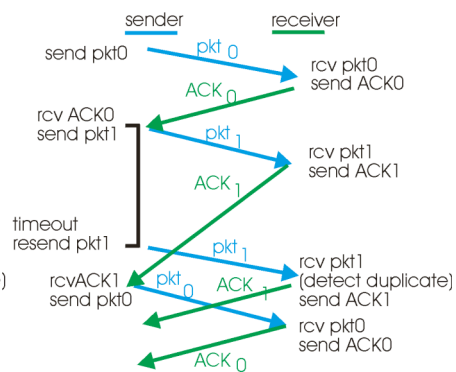
(b) lost packet

Transport Layer 3-39

rdt3.0 in action



(c) lost ACK



(d) premature timeout

Transport Layer 3-40

Performance of rdt3.0

r rdt3.0 works, but performance stinks

r example: 1 Gbps link, 15 ms e-e prop. delay, 1KB packet:

$$T_{\text{transmi}} = \frac{L \text{ (packet length in bits)}}{R \text{ (transmission rate, bps)}} = \frac{8\text{kb/pkt}}{10^{**9} \text{ b/sec}} = 8 \text{ microsec}$$

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

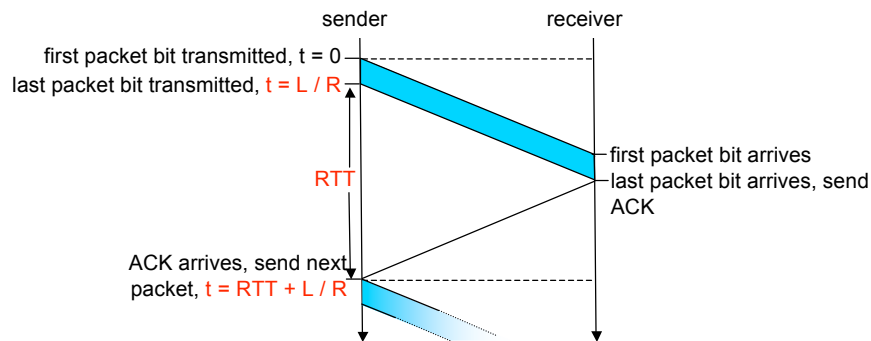
$$U_{\text{sender}} = \frac{L / R}{RTT + L / R} = \frac{.008}{30.008} = 0.00027$$

m 1KB pkt every 30 msec -> 33kB/sec thruput over 1 Gbps link

m network protocol limits use of physical resources!

Transport Layer 3-41

rdt3.0: stop-and-wait operation



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

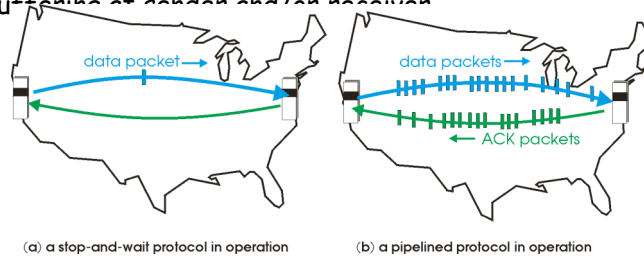
Transport Layer 3-42

Pipelined protocols

Pipelining: sender allows multiple, "in-flight", yet-to-be-acknowledged pkts

m range of sequence numbers must be increased

m buffering at sender and/or receiver



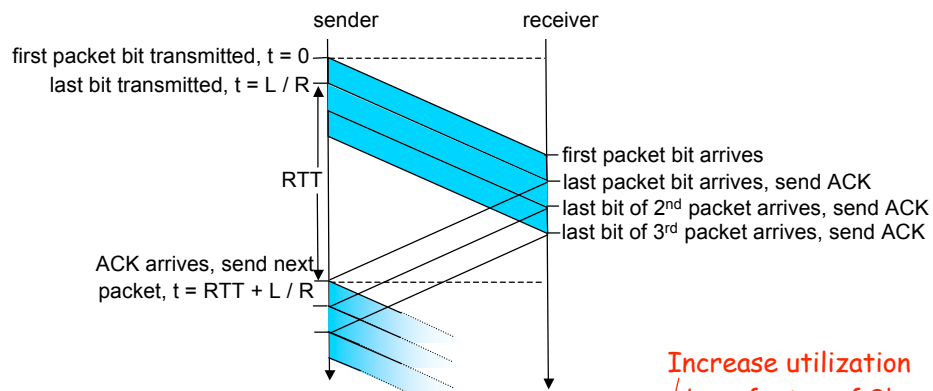
(a) a stop-and-wait protocol in operation

(b) a pipelined protocol in operation

r Two generic forms of pipelined protocols: *go-Back-N*, *selective repeat*

Transport Layer 3-43

Pipelining: increased utilization



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

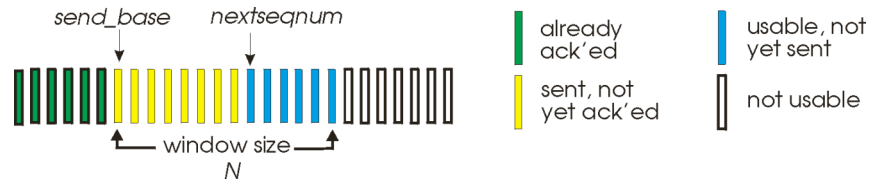
Increase utilization
by a factor of 3!

Transport Layer 3-44

Go-Back-N

Sender:

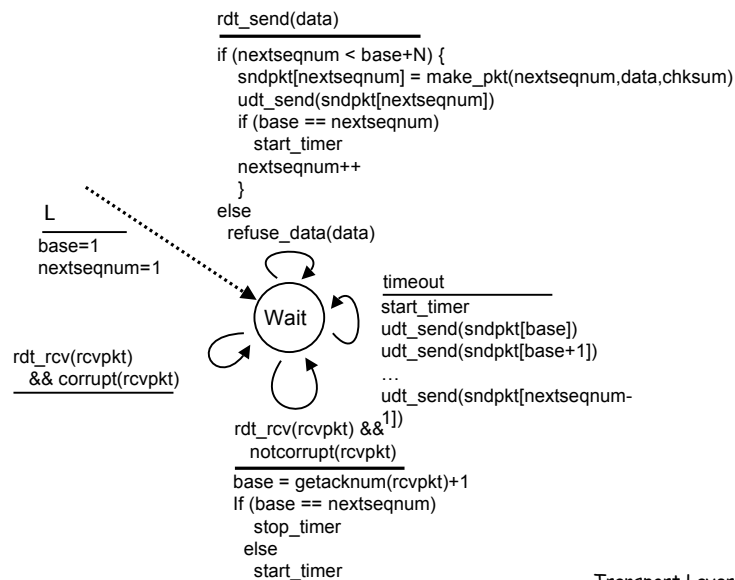
- r k-bit seq # in pkt header
- r "window" of up to N, consecutive unack'ed pkts allowed



- r ACK(n): ACKs all pkts up to, including seq # n - "cumulative ACK"
- m may receive duplicate ACKs (see receiver)
- r timer for each in-flight pkt
- r timeout(n): retransmit pkt n and all higher seq # pkts in window

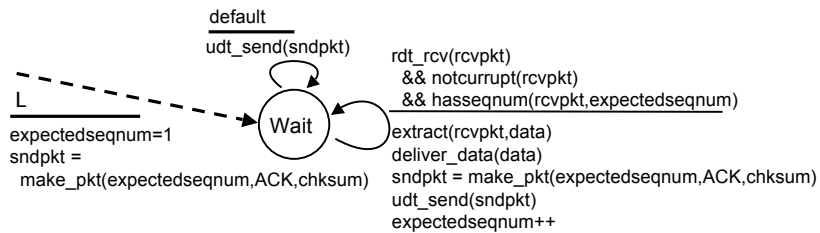
Transport Layer 3-45

GBN: sender extended FSM



Transport Layer 3-46

GBN: receiver extended FSM



ACK-only: always send ACK for correctly-received pkt with highest *in-order* seq #

m may generate duplicate ACKs

m need only remember **expectedseqnum**

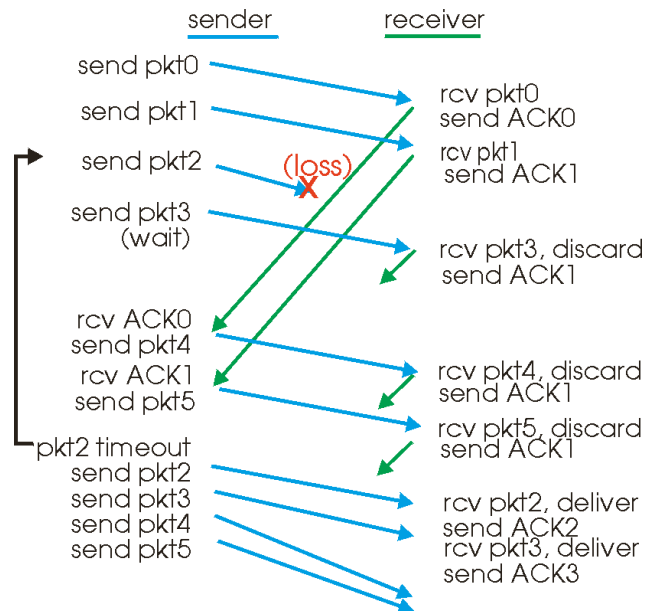
r out-of-order pkt:

m discard (don't buffer) -> **no receiver buffering!**

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

Transport Layer 3-47

GBN in action



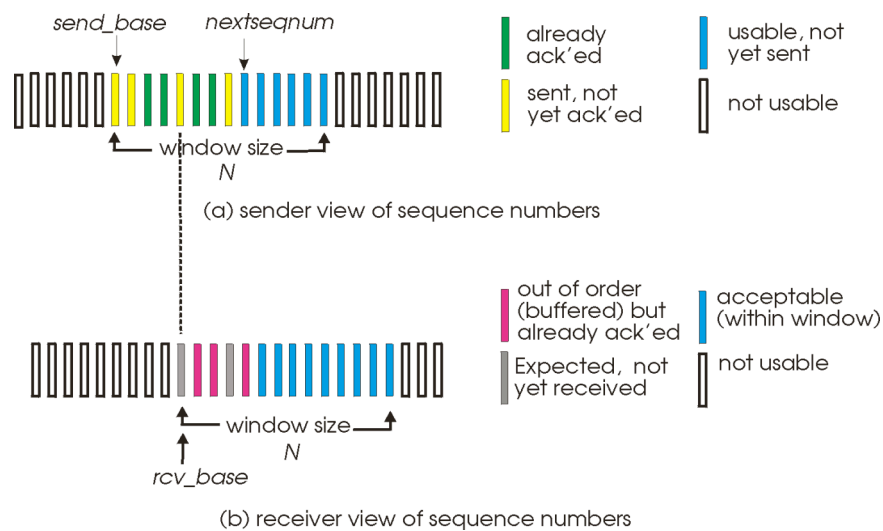
Transport Layer 3-48

Selective Repeat

- r receiver *individually* acknowledges all correctly received pkts
 - m buffers pkts, as needed, for eventual in-order delivery to upper layer
- r sender only resends pkts for which ACK not received
 - m sender timer for each unACKed pkt
- r sender window
 - m N consecutive seq #'s
 - m again limits seq #'s of sent, unACKed pkts

Transport Layer 3-49

Selective repeat: sender, receiver windows



Transport Layer 3-50

Selective repeat

sender

data from above :

- r if next available seq # in window, send pkt

timeout(n):

- r resend pkt n, restart timer

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

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

receiver

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

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

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

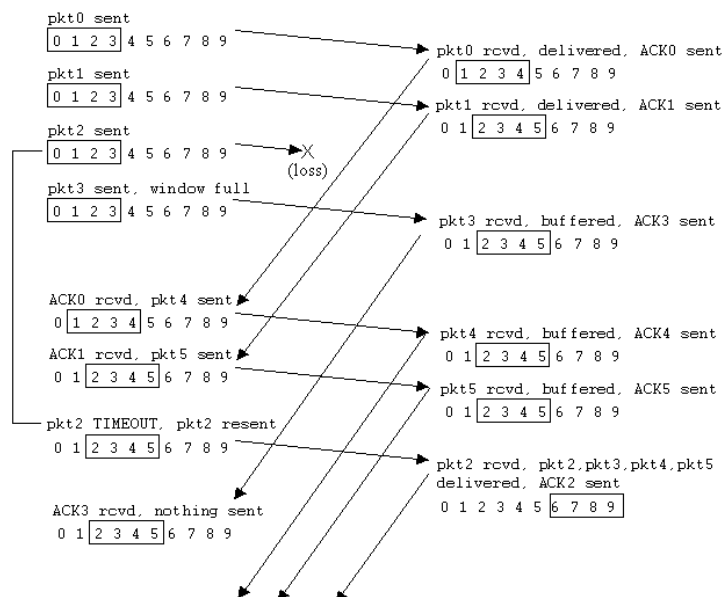
- r ACK(n)

otherwise:

- r ignore

Transport Layer 3-51

Selective repeat in action



† Layer 3-52

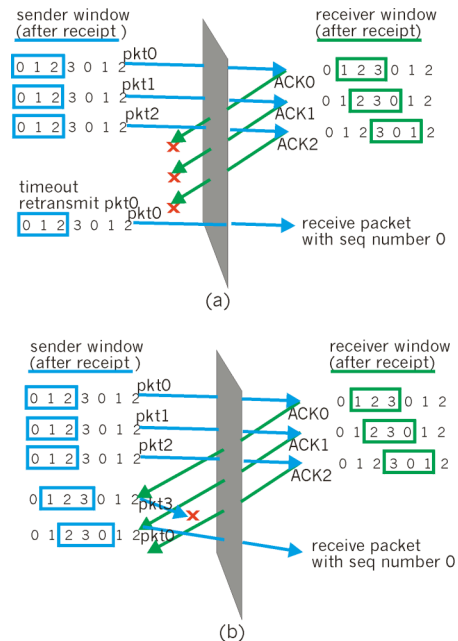
Selective repeat: dilemma

Example:

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

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

Q what relationship between seq # size and window size?



Transport Layer 3-53

Chapter 3 outline

- r 3.1 Transport-layer services
- r 3.2 Multiplexing and demultiplexing
- r 3.3 Connectionless transport: UDP
- r 3.4 Principles of reliable data transfer
- r 3.5 Connection-oriented transport: TCP
 - m segment structure
 - m reliable data transfer
 - m flow control
 - m connection management
- r 3.6 Principles of congestion control
- r 3.7 TCP congestion control

Transport Layer 3-54

TCP: Overview

RFCs: 793, 1122, 1323, 2018, 2581

r point-to-point:

- m one sender, one receiver

r reliable, in-order byte stream:

- m no "message boundaries"

r pipelined:

- m TCP congestion and flow control set window size

r send & receive buffers



r full duplex data:

- m bi-directional data flow in same connection

- m MSS: maximum segment size

r connection-oriented:

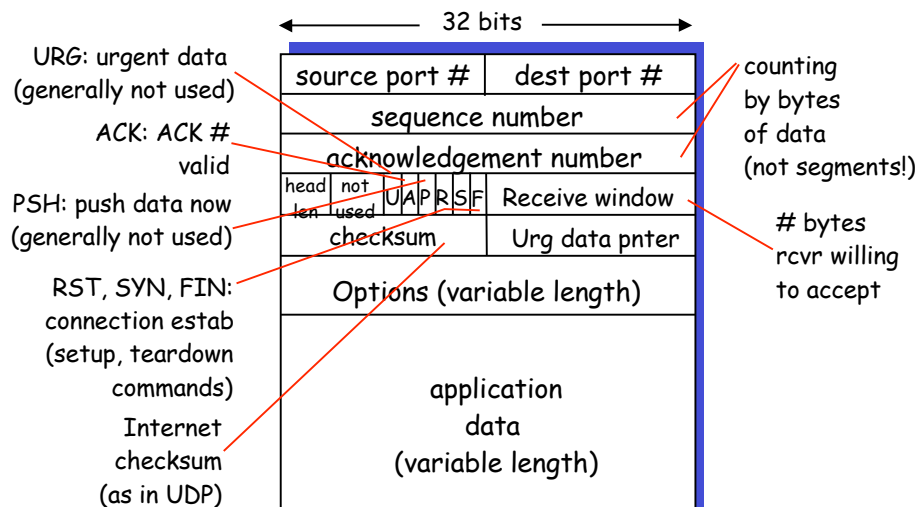
- m handshaking (exchange of control msgs) init's sender, receiver state before data exchange

r flow controlled:

- m sender will not overwhelm receiver

Transport Layer 3-55

TCP segment structure



Transport Layer 3-56

TCP seq. #'s and ACKs

Seq. #'s:

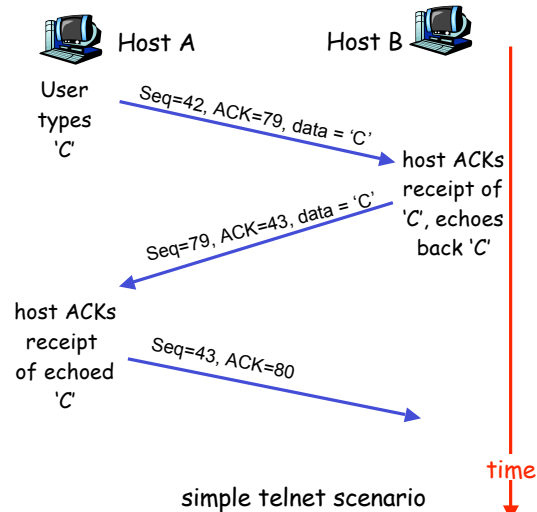
- m byte stream
"number" of first
byte in segment's
data

ACKs:

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

Q: how receiver handles
out-of-order segments

A: TCP spec doesn't say, -
up to implementor



Transport Layer 3-57

TCP Round Trip Time and Timeout

Q: how to set TCP
timeout value?

- r longer than RTT
 - m but RTT varies
- r too short: premature
timeout
 - m unnecessary
retransmissions
- r too long: slow reaction
to segment loss

Q: how to estimate RTT?

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

Transport Layer 3-58

TCP Round Trip Time and Timeout

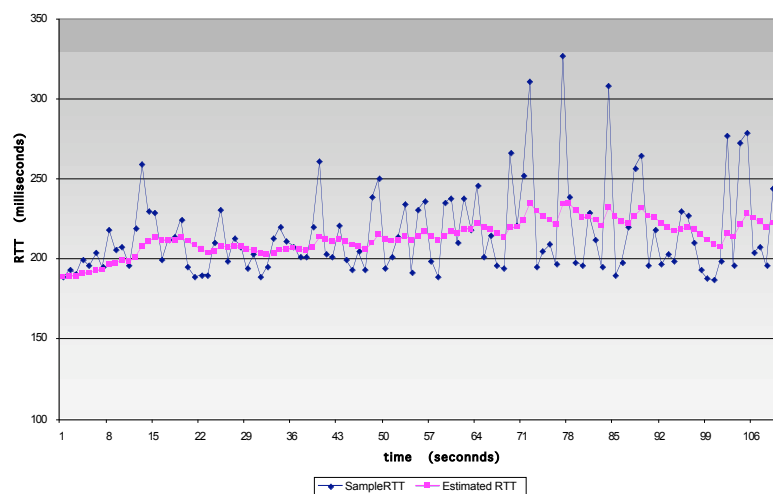
$$\text{EstimatedRTT} = (1 - \alpha) * \text{EstimatedRTT} + \alpha * \text{SampleRTT}$$

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

Transport Layer 3-59

Example RTT estimation:

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



Transport Layer 3-60

TCP Round Trip Time and Timeout

Setting the timeout

- r EstimatedRTT plus "safety margin"
 - m large variation in EstimatedRTT -> larger safety margin
- r first estimate of how much SampleRTT deviates from EstimatedRTT:
$$\text{DevRTT} = (1-\beta) * \text{DevRTT} + \beta * |\text{SampleRTT} - \text{EstimatedRTT}|$$

(typically, $\beta = 0.25$)

Then set timeout interval:

`TimeoutInterval = EstimatedRTT + 4*DevRTT`

Transport Layer 3-61

Chapter 3 outline

- | | |
|---|---|
| <ul style="list-style-type: none">r 3.1 Transport-layer servicesr 3.2 Multiplexing and demultiplexingr 3.3 Connectionless transport: UDPr 3.4 Principles of reliable data transfer | <ul style="list-style-type: none">r 3.5 Connection-oriented transport: TCP<ul style="list-style-type: none">m segment structurem reliable data transferm flow controlm connection managementr 3.6 Principles of congestion controlr 3.7 TCP congestion control |
|---|---|

Transport Layer 3-62

TCP reliable data transfer

- r TCP creates rdt service on top of IP's unreliable service
- r Pipelined segments
- r Cumulative acks
- r TCP uses single retransmission timer
- r Retransmissions are triggered by:
 - m timeout events
 - m duplicate acks
- r Initially consider simplified TCP sender:
 - m ignore duplicate acks
 - m ignore flow control, congestion control

Transport Layer 3-63

TCP sender events:

data rcvd from app:

- r Create segment with seq #
- r seq # is byte-stream number of first data byte in segment
- r start timer if not already running (think of timer as for oldest unacked segment)
- r expiration interval: TimeoutInterval

timeout:

- r retransmit segment that caused timeout
- r restart timer

Ack rcvd:

- r If acknowledges previously unacked segments
 - m update what is known to be acked
 - m start timer if there are outstanding segments

Transport Layer 3-64


```

NextSeqNum = InitialSeqNum
SendBase = InitialSeqNum

```

```

loop (forever) {
  switch(event)

```

```

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

```

```

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

```

```

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

```

```

} /* end of loop forever */

```

TCP sender (simplified)

Comment:

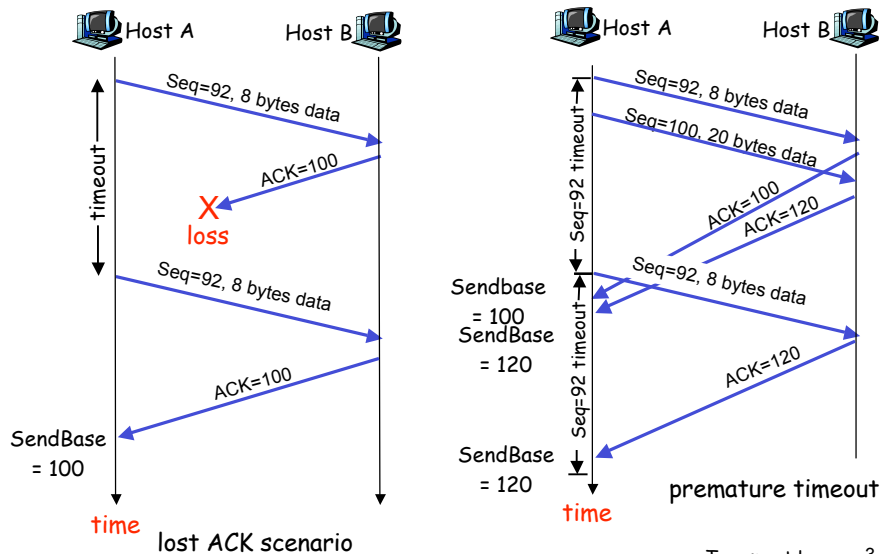
- SendBase-1: last cumulatively ack'ed byte

Example:

- SendBase-1 = 71; y = 73, so the rcvr wants 73+ ; y > SendBase, so that new data is acked

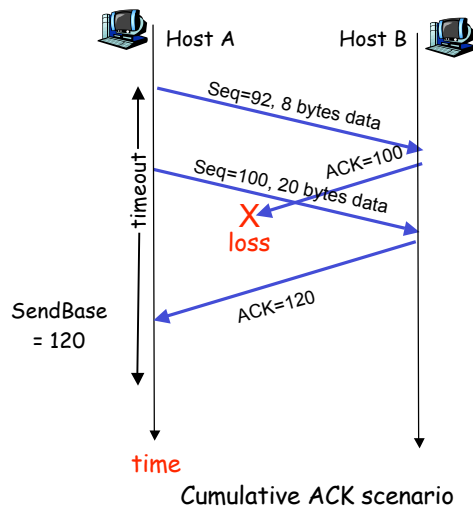
Transport Layer 3-65

TCP: retransmission scenarios



Transport Layer 3-66

TCP retransmission scenarios (more)



Transport Layer 3-67

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-expected seq. # . Gap detected	Immediately send <i>duplicate ACK</i> , 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

Transport Layer 3-68

Fast Retransmit

- r Time-out period often relatively long:
 - m long delay before resending lost packet
- r Detect lost segments via duplicate ACKs.
 - m Sender often sends many segments back-to-back
 - m If segment is lost, there will likely be many duplicate ACKs.
- r If sender receives 3 ACKs for the same data, it supposes that segment after ACKed data was lost:
 - m fast retransmit: resend segment before timer expires

Transport Layer 3-69

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

Transport Layer 3-70

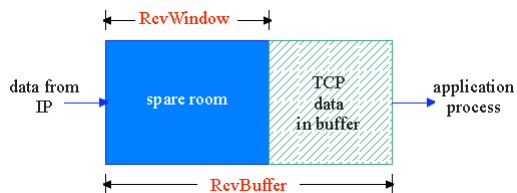
Chapter 3 outline

- r 3.1 Transport-layer services
- r 3.2 Multiplexing and demultiplexing
- r 3.3 Connectionless transport: UDP
- r 3.4 Principles of reliable data transfer
- r 3.5 Connection-oriented transport: TCP
 - m segment structure
 - m reliable data transfer
 - m **flow control**
 - m connection management
- r 3.6 Principles of congestion control
- r 3.7 TCP congestion control

Transport Layer 3-71

TCP Flow Control

- r receive side of TCP connection has a receive buffer:



- r app process may be slow at reading from buffer

flow control

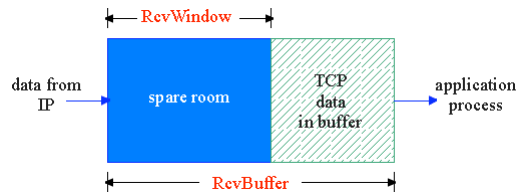
sender won't overflow receiver's buffer by transmitting too

much,
too fast

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

Transport Layer 3-72

TCP Flow control: how it works



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

- r spare room in buffer
- = `RcvWindow`
- = `RcvBuffer - [LastByteRcvd - LastByteRead]`

- r Rcvr advertises spare room by including value of `RcvWindow` in segments
- r Sender limits unACKed data to `RcvWindow`
 - m guarantees receive buffer doesn't overflow

Transport Layer 3-73

Chapter 3 outline

- r 3.1 Transport-layer services
- r 3.2 Multiplexing and demultiplexing
- r 3.3 Connectionless transport: UDP
- r 3.4 Principles of reliable data transfer
- r 3.5 Connection-oriented transport: TCP
 - m segment structure
 - m reliable data transfer
 - m flow control
 - m connection management
- r 3.6 Principles of congestion control
- r 3.7 TCP congestion control

Transport Layer 3-74

TCP Connection Management

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

r initialize TCP variables:

m seq. #s

m buffers, flow control info (e.g. RcvWindow)

r *client*: connection initiator

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

r *server*: contacted by client

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

Three way handshake:

Step 1: client host sends TCP SYN segment to server

m specifies initial seq #

m no data

Step 2: server host receives SYN, replies with SYNACK segment

m server allocates buffers

m specifies server initial seq. #

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

Transport Layer 3-75

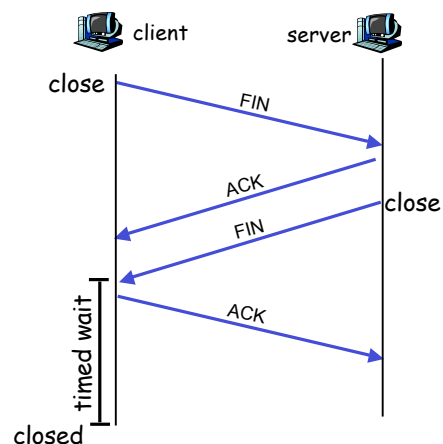
TCP Connection Management (cont.)

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.



Transport Layer 3-76

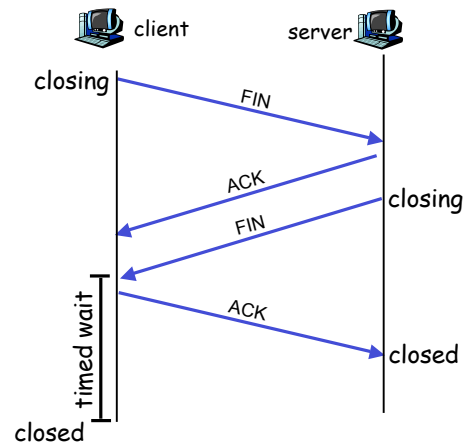
TCP Connection Management (cont.)

Step 3: client receives FIN, replies with ACK.

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

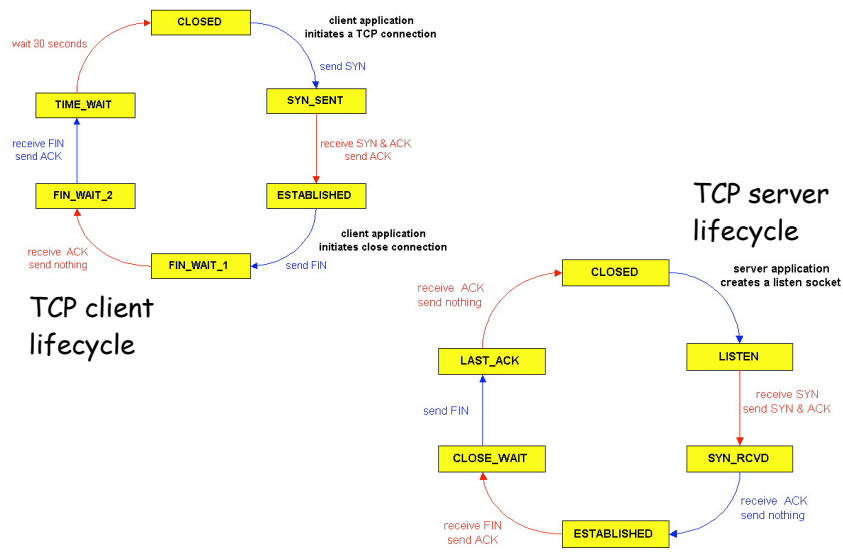
Step 4: server, receives ACK. Connection closed.

Note: with small modification, can handle simultaneous FINs.



Transport Layer 3-77

TCP Connection Management (cont)



Transport Layer 3-78

Chapter 3 outline

- r 3.1 Transport-layer services
- r 3.2 Multiplexing and demultiplexing
- r 3.3 Connectionless transport: UDP
- r 3.4 Principles of reliable data transfer
- r 3.5 Connection-oriented transport: TCP
 - m segment structure
 - m reliable data transfer
 - m flow control
 - m connection management
- r 3.6 Principles of congestion control
- r 3.7 TCP congestion control

Transport Layer 3-79

Principles of Congestion Control

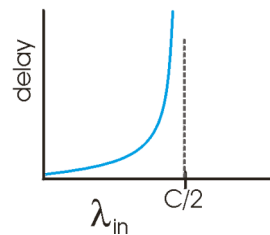
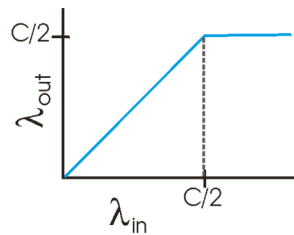
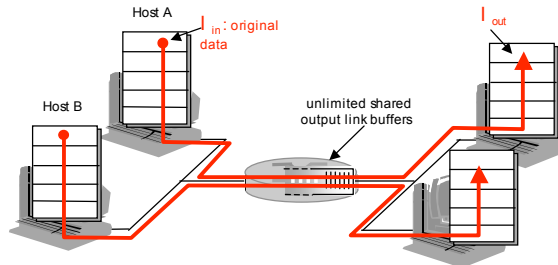
Congestion:

- r informally: "too many sources sending too much data too fast for *network* to handle"
- r different from flow control!
- r manifestations:
 - m lost packets (buffer overflow at routers)
 - m long delays (queueing in router buffers)
- r a top-10 problem!

Transport Layer 3-80

Causes/costs of congestion: scenario 1

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

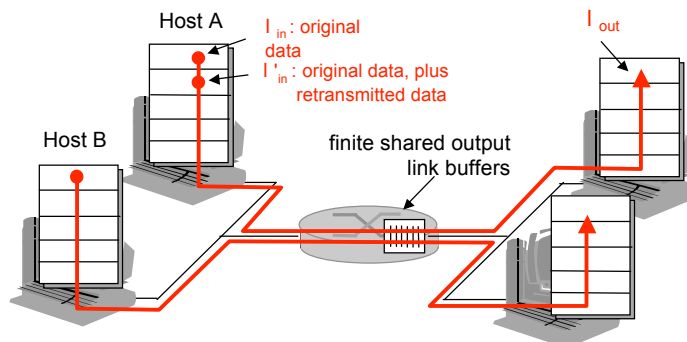


- large delays when congested
- maximum achievable throughput

Transport Layer 3-81

Causes/costs of congestion: scenario 2

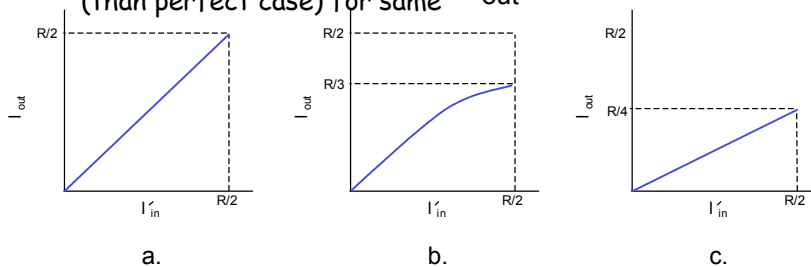
- one router, *finite* buffers
- sender retransmission of lost packet



Transport Layer 3-82

Causes/costs of congestion: scenario 2

- always: $I_{in} = I_{out}$ (goodput)
- "perfect" retransmission only when loss: $I'_{in} > I_{out}$,
- retransmission of delayed (not lost) packet makes I_{in} larger (than perfect case) for same I_{out}



"costs" of congestion:

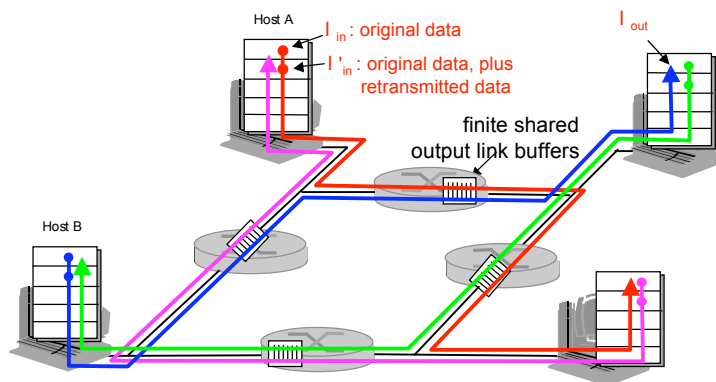
- more work (retrans) for given "goodput"
- unnneeded retransmissions: link carries multiple copies of pkt

Transport Layer 3-83

Causes/costs of congestion: scenario 3

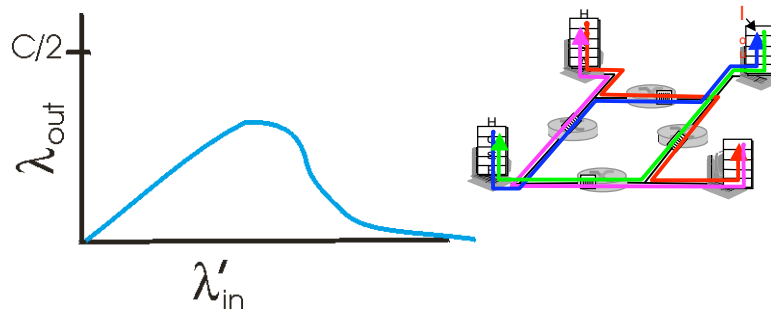
- four senders
- multihop paths
- timeout/retransmit

Q: what happens as I_{in} and I'_{in} increase ?



Transport Layer 3-84

Causes/costs of congestion: scenario 3



Another "cost" of congestion:

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

Transport Layer 3-85

Approaches towards congestion control

Two broad approaches towards congestion control:

End-end congestion control:

- r no explicit feedback from network
- r congestion inferred from end-system observed loss, delay
- r approach taken by TCP

Network-assisted congestion control:

- r routers provide feedback to end systems
 - m single bit indicating congestion (SNA, DECbit, TCP/IP ECN, ATM)
 - m explicit rate sender should send at

Transport Layer 3-86

Case study: ATM ABR congestion control

ABR: available bit rate:

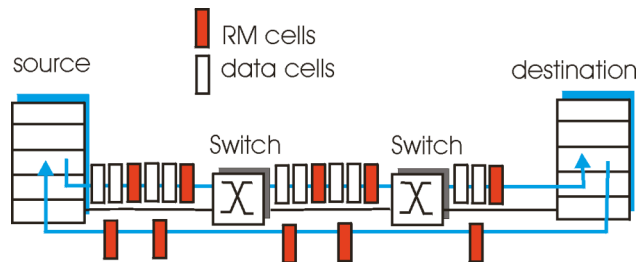
- r "elastic service"
- r if sender's path "underloaded":
 - m sender should use available bandwidth
- r if sender's path congested:
 - m sender throttled to minimum guaranteed rate

RM (resource management) cells:

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

Transport Layer 3-87

Case study: ATM ABR congestion control



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

Transport Layer 3-88

Chapter 3 outline

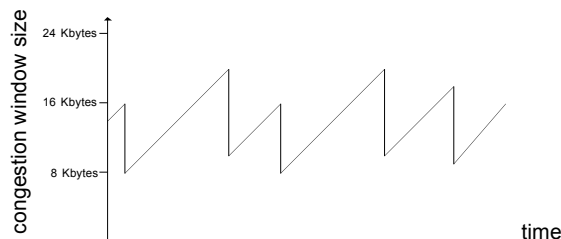
- r 3.1 Transport-layer services
- r 3.2 Multiplexing and demultiplexing
- r 3.3 Connectionless transport: UDP
- r 3.4 Principles of reliable data transfer
- r 3.5 Connection-oriented transport: TCP
 - m segment structure
 - m reliable data transfer
 - m flow control
 - m connection management
- r 3.6 Principles of congestion control
- r 3.7 TCP congestion control

Transport Layer 3-89

TCP congestion control: additive increase, multiplicative decrease

- r **Approach:** increase transmission rate (window size), probing for usable bandwidth, until loss occurs
 - m **additive increase:** increase **CongWin** by 1 MSS every RTT until loss detected
 - m **multiplicative decrease:** cut **CongWin** in half after loss

Saw tooth behavior: probing for bandwidth



Transport Layer 3-90

TCP Congestion Control: details

- r sender limits transmission:

$$\text{LastByteSent} - \text{LastByteAcked} \leq \text{CongWin}$$
 - r Roughly,

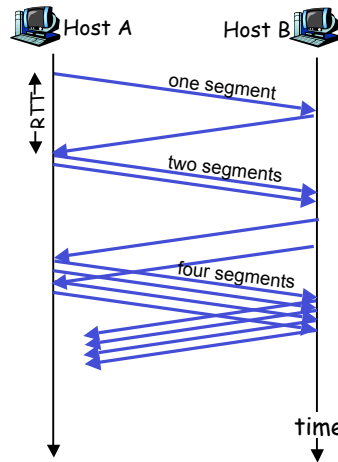
$$\text{rate} = \frac{\text{CongWin}}{\text{RTT}} \text{ Bytes/sec}$$
 - r CongWin is dynamic, function of perceived network congestion
- How does sender perceive congestion?
- r loss event = timeout or 3 duplicate acks
 - r TCP sender reduces rate (CongWin) after loss event
- three mechanisms:
- m AIMD
 - m slow start
 - m conservative after timeout events
- Transport Layer 3-91

TCP Slow Start

- r When connection begins, CongWin = 1 MSS
 - m Example: MSS = 500 bytes & RTT = 200 msec
 - m initial rate = 20 kbps
- r available bandwidth may be \gg MSS/RTT
 - m desirable to quickly ramp up to respectable rate
- r When connection begins, increase rate exponentially fast until first loss event

TCP Slow Start (more)

- r When connection begins, increase rate exponentially until first loss event:
 - m double CongWin every RTT
 - m done by incrementing CongWin for every ACK received
- r Summary: initial rate is slow but ramps up exponentially fast



Transport Layer 3-93

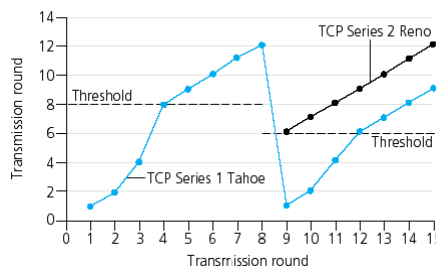
Refinement

Q: When should the exponential increase switch to linear?

A: When CongWin gets to 1/2 of its value before timeout.

Implementation:

- r Variable Threshold
- r At loss event, Threshold is set to 1/2 of CongWin just before loss event



Transport Layer 3-94

Refinement: inferring loss

- r After 3 dup ACKs:
 - m CongWin is cut in half
 - m window then grows linearly
- r But after timeout event:
 - m CongWin instead set to 1 MSS;
 - m window then grows exponentially
 - m to a threshold, then grows linearly

Philosophy:

- 3 dup ACKs indicates network capable of delivering some segments
- timeout indicates a "more alarming" congestion scenario

Transport Layer 3-95

Summary: TCP Congestion Control

- r When CongWin is below Threshold, sender in **slow-start** phase, window grows exponentially.
- r When CongWin is above Threshold, sender is in **congestion-avoidance** phase, window grows linearly.
- r When a **triple duplicate ACK** occurs, Threshold set to CongWin/2 and CongWin set to Threshold.
- r When **timeout** occurs, Threshold set to CongWin/2 and CongWin is set to 1 MSS.

Transport Layer 3-96

TCP sender congestion control

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

Transport Layer 3-97

TCP throughput

- r What's the average throughput of TCP as a function of window size and RTT?
 - m Ignore slow start
- r Let W be the window size when loss occurs.
- r When window is W , throughput is W/RTT
- r Just after loss, window drops to $W/2$, throughput to $W/2\text{RTT}$.
- r Average throughput: $.75 W/\text{RTT}$

Transport Layer 3-98

TCP Futures: TCP over "long, fat pipes"

- r Example: 1500 byte segments, 100ms RTT, want 10 Gbps throughput
- r Requires window size $W = 83,333$ in-flight segments
- r Throughput in terms of loss rate:

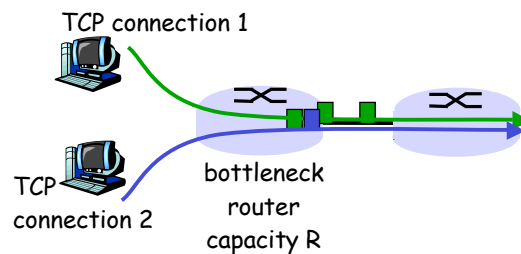
$$\frac{1.22 \cdot MSS}{RTT\sqrt{L}}$$

- r $\rightarrow L = 2 \cdot 10^{-10}$ **Wow**
- r New versions of TCP for high-speed needed!

Transport Layer 3-99

TCP Fairness

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

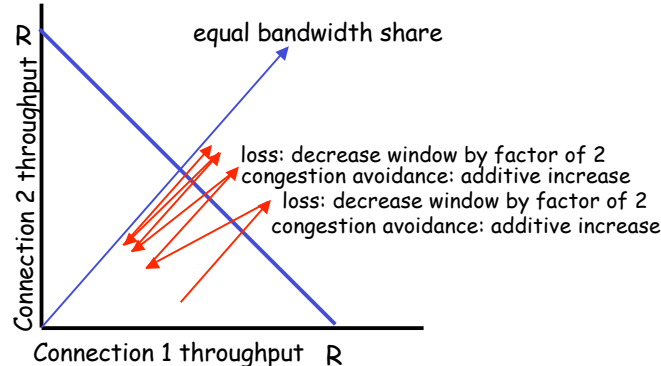


Transport Layer 3-100

Why is TCP fair?

Two competing sessions:

- r Additive increase gives slope of 1, as throughput increases
- r multiplicative decrease decreases throughput proportionally



Transport Layer 3-101

Fairness (more)

Fairness and UDP

- r Multimedia apps often do not use TCP
 - m do not want rate throttled by congestion control
- r Instead use UDP:
 - m pump audio/video at constant rate, tolerate packet loss
- r Research area: TCP friendly

Fairness and parallel TCP connections

- r nothing prevents app from opening parallel connections between 2 hosts.
- r Web browsers do this
- r Example: link of rate R supporting 9 cncctions;
 - m new app asks for 1 TCP, gets rate $R/10$
 - m new app asks for 11 TCPs, gets $R/2$!

Transport Layer 3-102

Delay modeling

Q: How long does it take to receive an object from a Web server after sending a request?

Ignoring congestion, delay is influenced by:

- r TCP connection establishment
- r data transmission delay
- r slow start

Notation, assumptions:

- r Assume one link between client and server of rate R
- r S : MSS (bits)
- r O : object size (bits)
- r no retransmissions (no loss, no corruption)

Window size:

- r First assume: fixed congestion window, W segments
- r Then dynamic window, modeling slow start

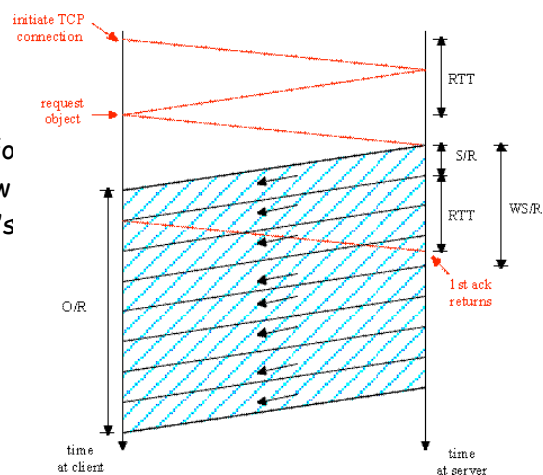
Transport Layer 3-103

Fixed congestion window (1)

First case:

$WS/R > RTT + S/R$: ACK for first segment in window returns before window's worth of data sent

$$\text{delay} = 2RTT + O/R$$



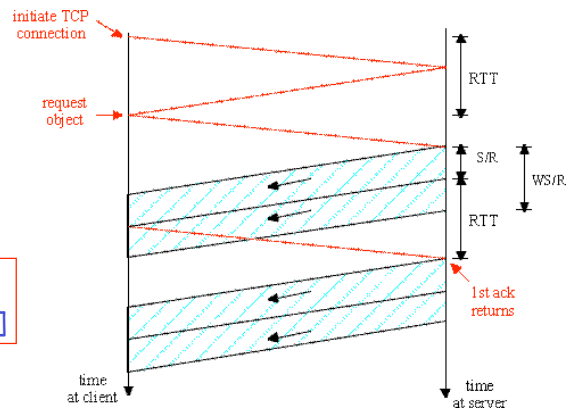
Transport Layer 3-104

Fixed congestion window (2)

Second case:

- $WS/R < RTT + S/R$: wait for ACK after sending window's worth of data sent

$$\text{delay} = 2RTT + O/R + (K-1)[S/R + RTT - WS/R]$$



Transport Layer 3-105

TCP Delay Modeling: Slow Start (1)

Now suppose window grows according to slow start

Will show that the delay for one object is:

$$\text{Latency} = 2RTT + \frac{O}{R} + P \left[RTT + \frac{S}{R} \right] - (2^P - 1) \frac{S}{R}$$

where P is the number of times TCP idles at server:

$$P = \min\{Q, K - 1\}$$

- where Q is the number of times the server idles if the object were of infinite size.
- and K is the number of windows that cover the object.

Transport Layer 3-106

TCP Delay Modeling: Slow Start (2)

Delay components:

- 2 RTT for connection estab and request
- O/R to transmit object
- time server idles due to slow start

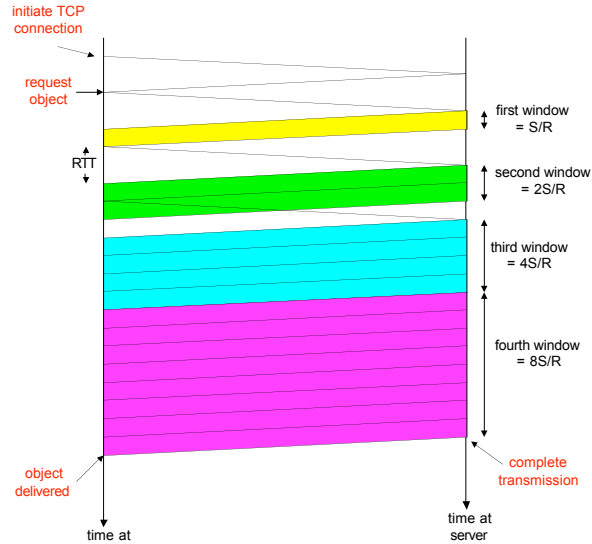
Server idles:

$$P = \min\{K-1, Q\} \text{ times}$$

Example:

- O/S = 15 segments
- K = 4 windows
- Q = 2
- P = min{K-1, Q} = 2

Server idles P=2 times



Transport Layer 3-107

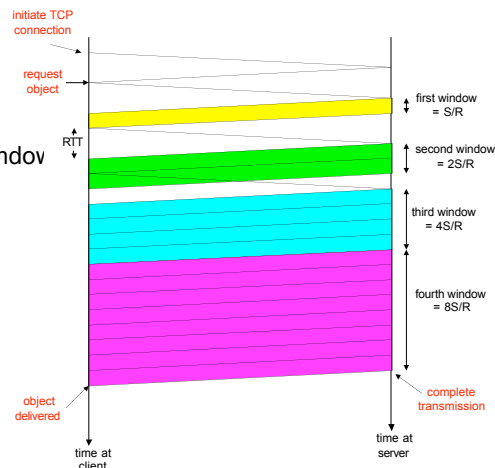
TCP Delay Modeling (3)

$\frac{S}{R} + RTT$ = time from where server starts to send segment until server receives acknowledgment

$2^{k-1} \frac{S}{R}$ = time to transmit the kth window

$\left[\frac{S}{R} + RTT - 2^{k-1} \frac{S}{R} \right]^+ = \text{idle time after the kth window}$

$$\begin{aligned} \text{delay} &= \frac{O}{R} + 2RTT + \sum_{p=1}^P \text{idleTime}_p \\ &= \frac{O}{R} + 2RTT + \sum_{k=1}^P \left[\frac{S}{R} + RTT - 2^{k-1} \frac{S}{R} \right] \\ &= \frac{O}{R} + 2RTT + P \left[RTT + \frac{S}{R} \right] - (2^P - 1) \frac{S}{R} \end{aligned}$$



Transport Layer 3-108

TCP Delay Modeling (4)

Recall K = number of windows that cover object

How do we calculate K ?

$$\begin{aligned} K &= \min\{k: 2^0 S + 2^1 S + L + 2^{k-1} S \geq O\} \\ &= \min\{k: 2^0 + 2^1 + L/S + 2^{k-1} \geq O/S\} \\ &= \min\{k: 2^k - 1 \geq \frac{O}{S}\} \\ &= \min\{k: k \geq \log_2(\frac{O}{S} + 1)\} \\ &= \left\lceil \log_2(\frac{O}{S} + 1) \right\rceil \end{aligned}$$

Calculation of Q , number of idles for infinite-size object, is similar (see HW).

Transport Layer 3-109

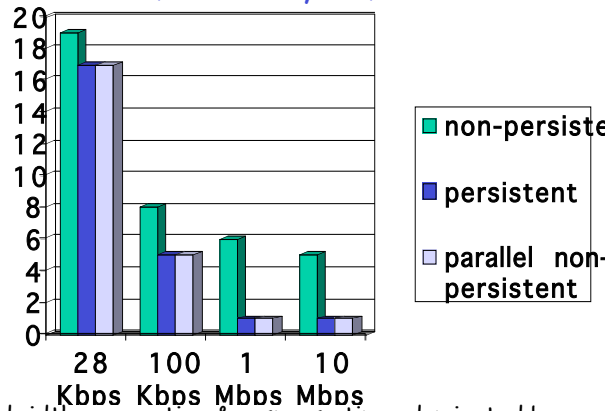
HTTP Modeling

- r Assume Web page consists of:
 - m 1 base HTML page (of size O bits)
 - m M images (each of size O bits)
- r Non-persistent HTTP:
 - m $M+1$ TCP connections in series
 - m Response time = $(M+1)O/R + (M+1)2RTT + \text{sum of idle times}$
- r Persistent HTTP:
 - m 2 RTT to request and receive base HTML file
 - m 1 RTT to request and receive M images
 - m Response time = $(M+1)O/R + 3RTT + \text{sum of idle times}$
- r Non-persistent HTTP with X parallel connections
 - m Suppose M/X integer.
 - m 1 TCP connection for base file
 - m M/X sets of parallel connections for images.
 - m Response time = $(M+1)O/R + (M/X + 1)2RTT + \text{sum of idle times}$

Transport Layer 3-110

HTTP Response time (in seconds)

RTT = 100 msec, O = 5 Kbytes, M=10 and X=5



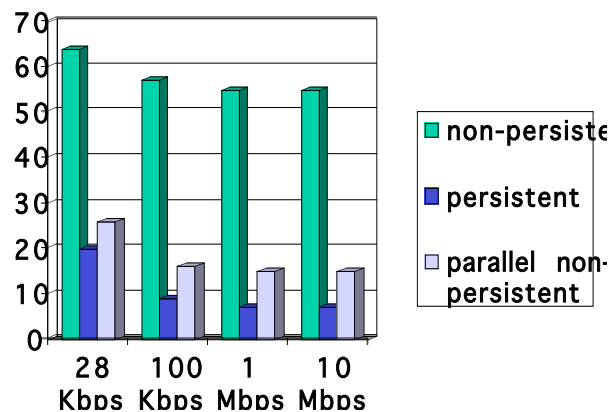
For low bandwidth, connection & response time dominated by transmission time.

Persistent connections only give minor improvement over parallel connections.

Transport Layer 3-111

HTTP Response time (in seconds)

RTT = 1 sec, O = 5 Kbytes, M=10 and X=5



For larger RTT, response time dominated by TCP establishment & slow start delays. Persistent connections now give important improvement: particularly in high delay•bandwidth networks.

Transport Layer 3-112

Chapter 3: Summary

- r principles behind transport layer services:
 - m multiplexing, demultiplexing
 - m reliable data transfer
 - m flow control
 - m congestion control
 - r instantiation and implementation in the Internet
 - m UDP
 - m TCP
- Next:
- r leaving the network "edge" (application, transport layers)
 - r into the network "core"

Transport Layer 3-113