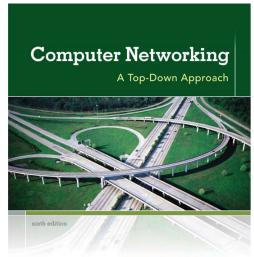
# Topic 6 Transport Layer



KUROSE ROSS

#### A note on the use of these ppt 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

©All material copyright 1996-2013 J.F Kurose and K.W. Ross, All Rights Reserved Computer
Networking: A Top
Down Approach
6<sup>th</sup> edition
Jim Kurose, Keith Ross
Addison-Wesley
March 2012

## Topic 6: 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

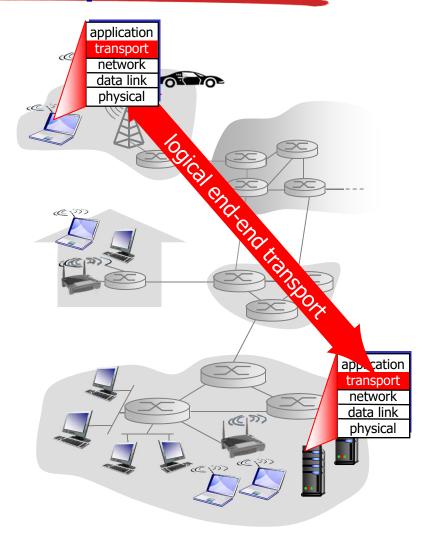
# **Outline**

- 6.1 transport-layer services
- 6.2 multiplexing and demultiplexing
- 6.3 connectionless transport: UDP
- 6.4 principles of reliable data transfer

- 6.5 connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- 6.6 principles of congestion control
- 6.7 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



# Transport vs. network layer

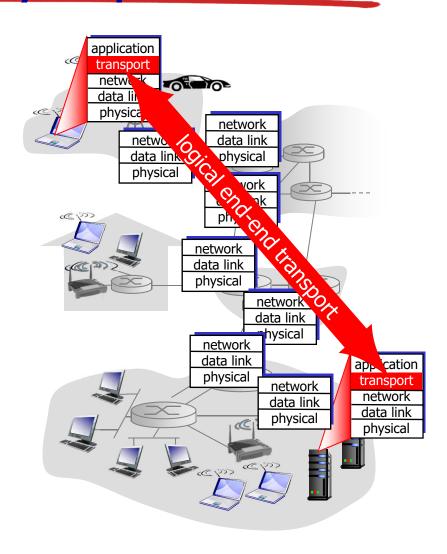
- network layer: logical communication between hosts
- transport layer: logical communication between processes
  - relies on, enhances, network layer services

#### household analogy:

- 12 kids in Ann's house sending letters to 12 kids in Bill's house:
- hosts = houses
- processes = kids
- app messages = letters in envelopes
- transport protocol = Ann and Bill who demux to inhouse 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



# **Outline**

- 6.1 transport-layer services
- 6.2 multiplexing and demultiplexing
- 6.3 connectionless transport: UDP
- 6.4 principles of reliable data transfer

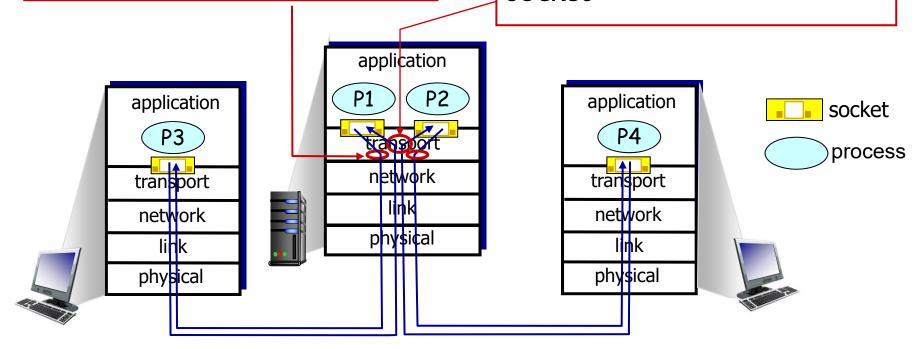
- 6.5 connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- 6.6 principles of congestion control
- 6.7 TCP congestion control

# Multiplexing/demultiplexing

#### multiplexing at sender:

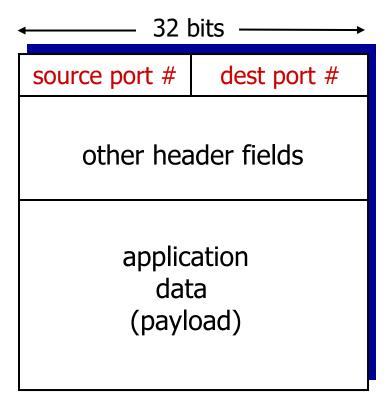
handle data from multiple sockets, add transport header (later used for demultiplexing) demultiplexing at receiver:

use header info to deliver received segments to correct socket



#### How demultiplexing works

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



TCP/UDP segment format

# Connectionless demultiplexing

created socket has host-local \* when creating datagram to port #:

DatagramSocket mySocket1 = new DatagramSocket(12534);

- send into UDP socket, must specify
  - destination IP address
  - destination port #

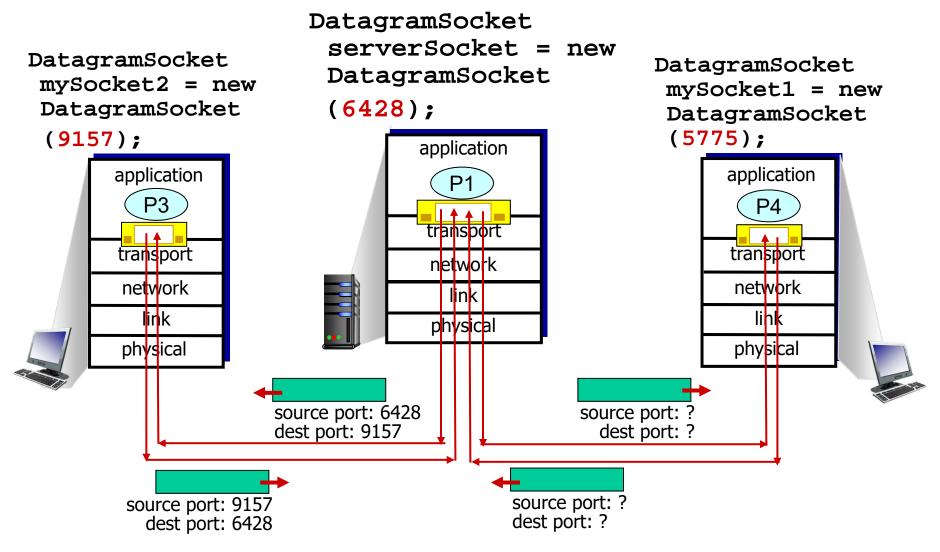
- when host receives UDP segment:
  - checks destination port # in segment



directs UDP segment to socket with that port #

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

## Connectionless demux: example

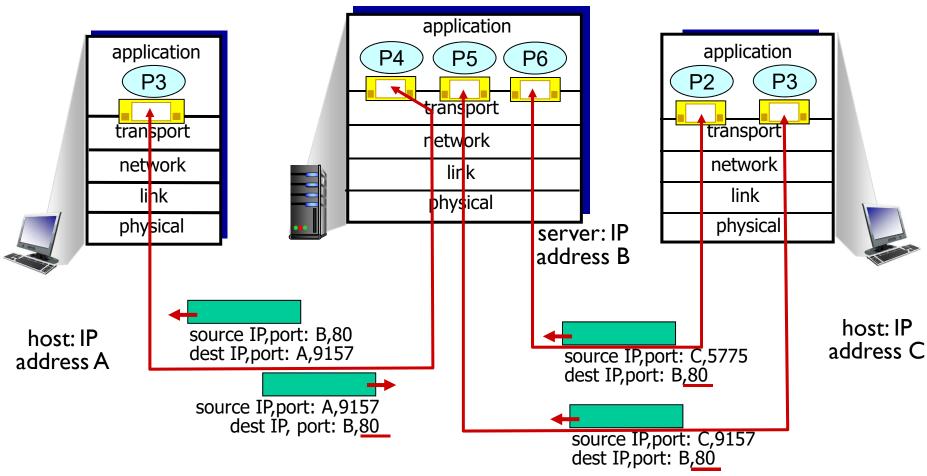


#### Connection-oriented demux

- 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

#### Connection-oriented demux: example



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

# <u>Outline</u>

- 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

# Applications and the underlying transport protocols

Application	Application-Layer Protocol	Underlying Transport Protocol
Remote terminal access	Telnet	TCP
Web	НТТР	TCP
File transfer	FTP	TCP
Remote file server	NFS	Typically UDP
Streaming multimedia	typically proprietory	UDP or TCP
Internet telephony	typically proprietory	UDP or TCP
Network Management	SNMP	Typically UDP
Routing Protocol	RIP	Typically UDP
Name translation	DNS	Typically UDP

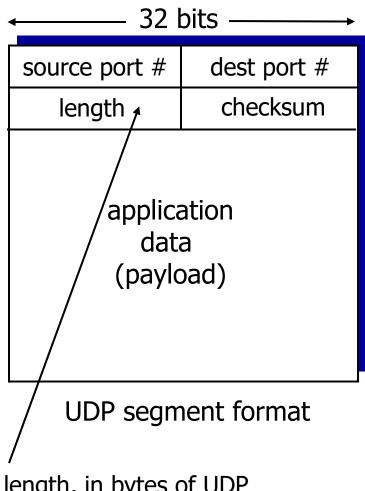
## UDP: User Datagram Protocol [RFC 768]

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

#### why is there a UDP?

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

## UDP: segment header



length, in bytes of UDP segment, including header

- 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** checksum

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

#### sender:

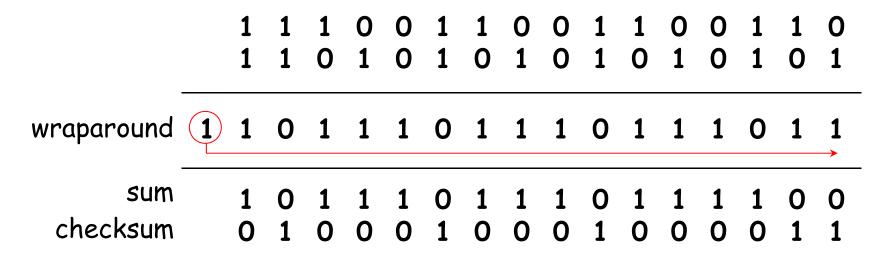
- treat segment contents, including header fields, as sequence of 16-bit integers
- checksum: addition (one's complement sum) of segment contents
- sender puts checksum value into UDP checksum field

#### receiver:

- compute checksum of received segment
- check if computed checksum equals checksum field value:
  - NO error detected
  - YES no error detected

## Internet checksum: example

example: add two 16-bit integers



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

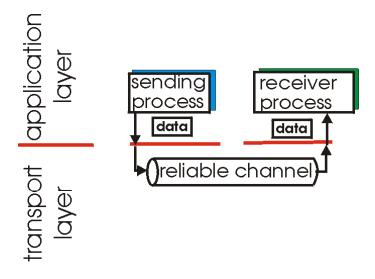
## **Outline**

- 6.1 transport-layer services
- 6.2 multiplexing and demultiplexing
- 6.3 connectionless transport: UDP
- 6.4 principles of reliable data transfer

- 6.5 connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- 6.6 principles of congestion control
- 6.7 TCP congestion control

## Principles of reliable data transfer

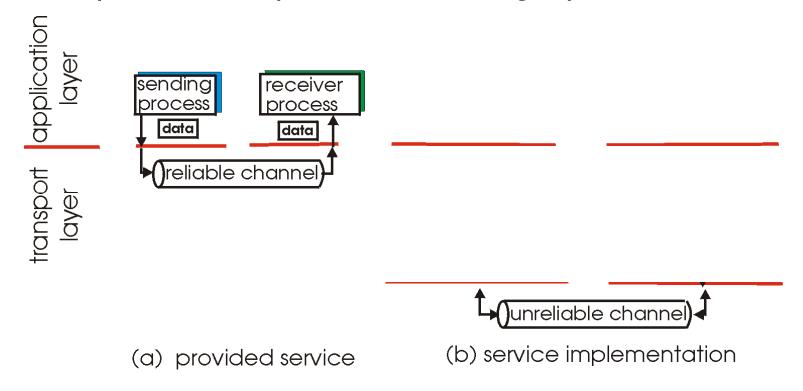
- important in application, transport, link layers
  - top-10 list of important networking topics!



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

## Principles of reliable data transfer

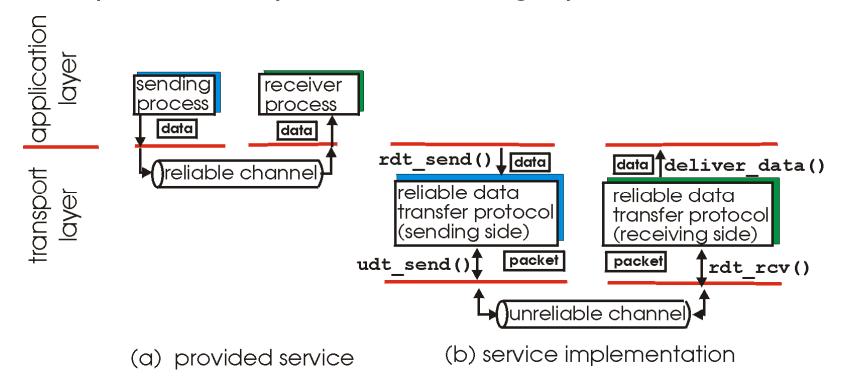
- important in application, transport, link layers
  - top-10 list of important networking topics!



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

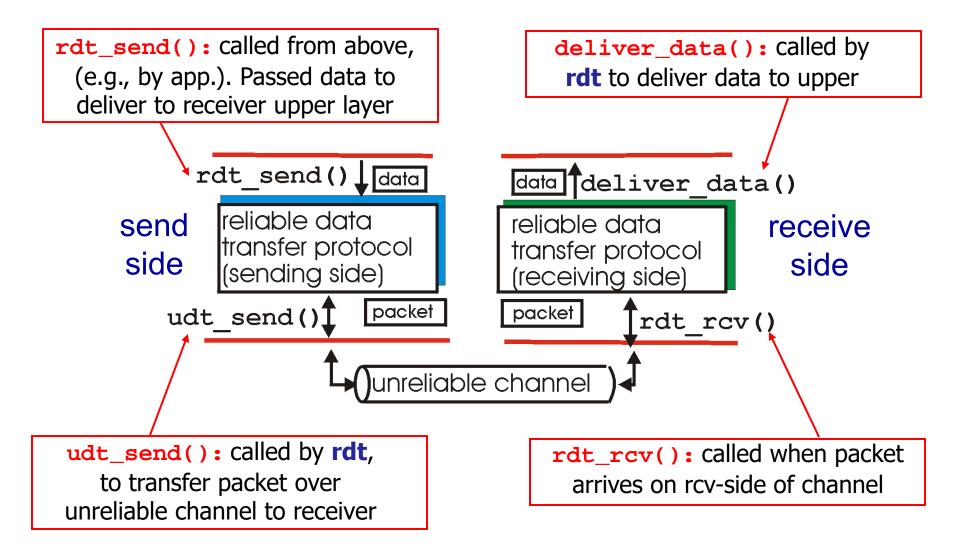
## Principles of reliable data transfer

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



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

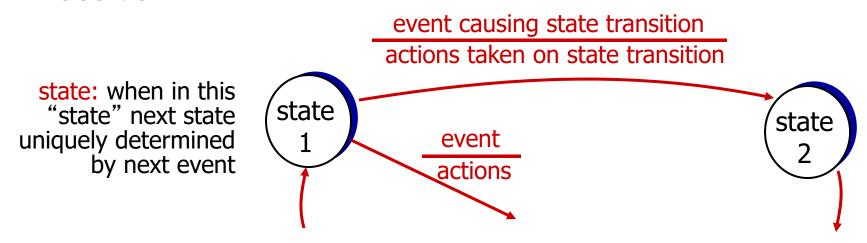
#### Reliable data transfer: getting started



#### Reliable data transfer: getting started

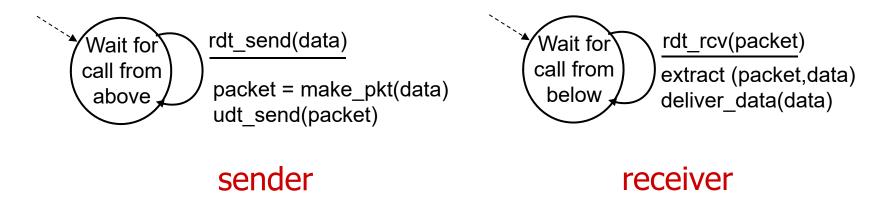
#### we'll:

- 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



#### rdt 1.0: 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



#### rdt2.0: channel with bit errors

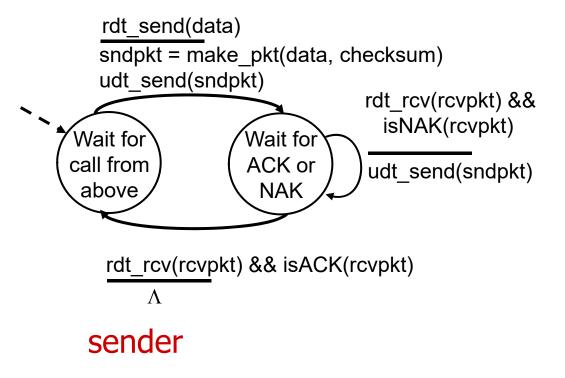
- underlying channel may flip bits in packet
  - Example: checksum to detect bit errors
- the question: how to recover from errors:

How do humans recover from "errors" during conversation?

#### rdt2.0: channel with bit errors

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

## rdt2.0: FSM specification



#### receiver

rdt\_rcv(rcvpkt) &&
corrupt(rcvpkt)

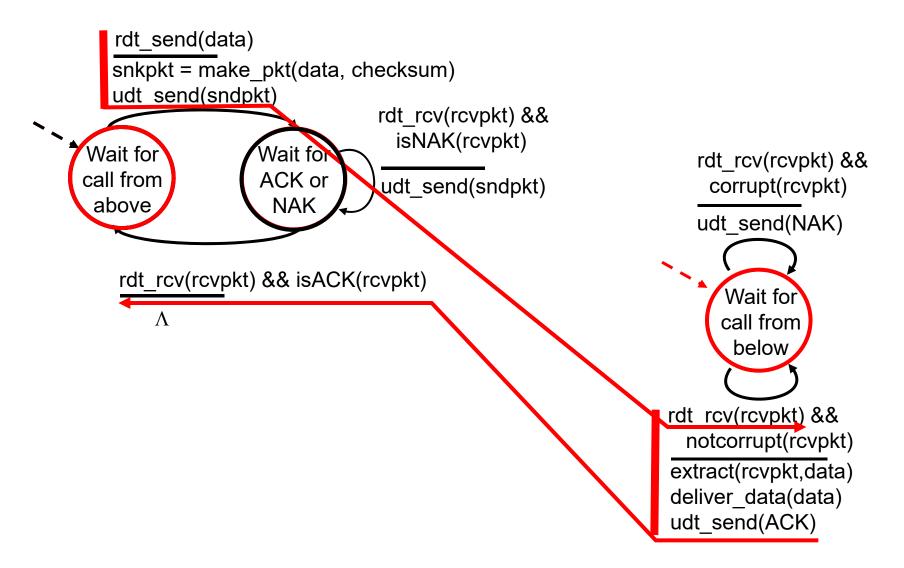
udt\_send(NAK)

Wait for
call from
below

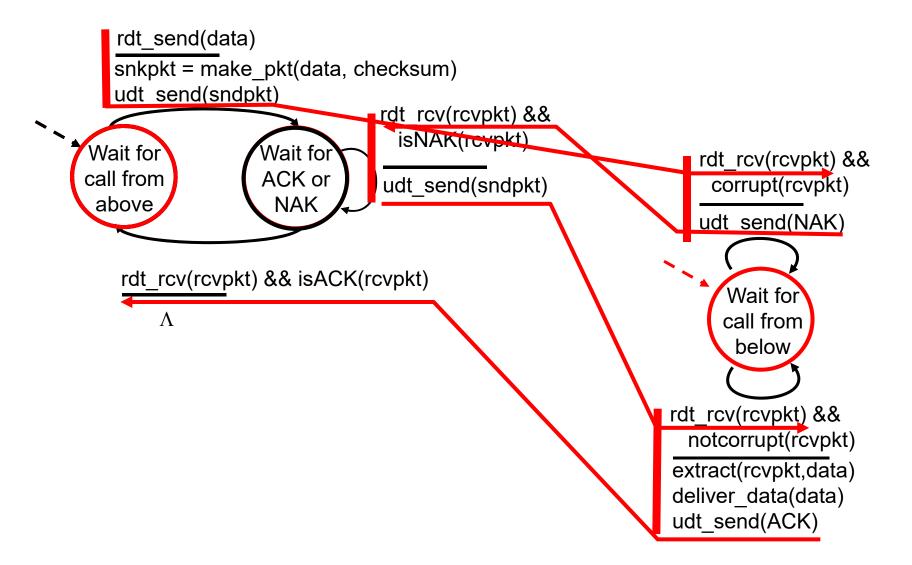
rdt\_rcv(rcvpkt) &&
notcorrupt(rcvpkt)

extract(rcvpkt,data)
deliver\_data(data)
udt\_send(ACK)

### rdt2.0: operation with no errors



#### rdt2.0: error scenario



### rdt2.0 has a fatal flaw!

# what happens if ACK/NAK corrupted?

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

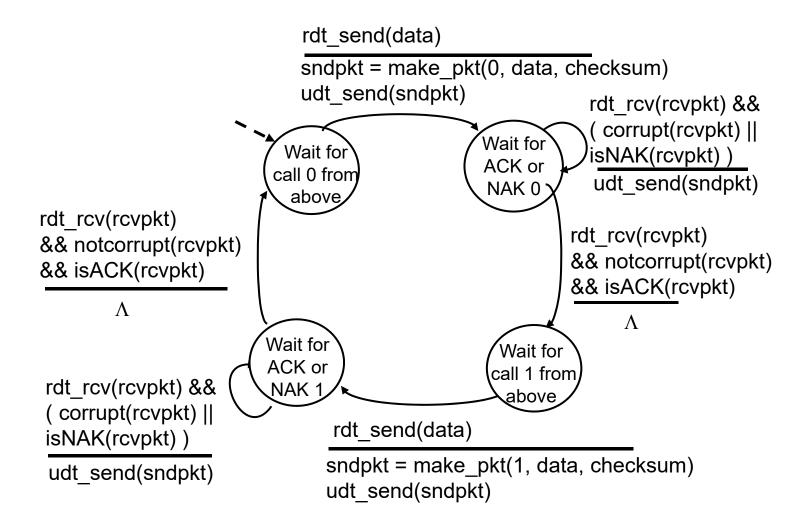
#### handling duplicates:

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

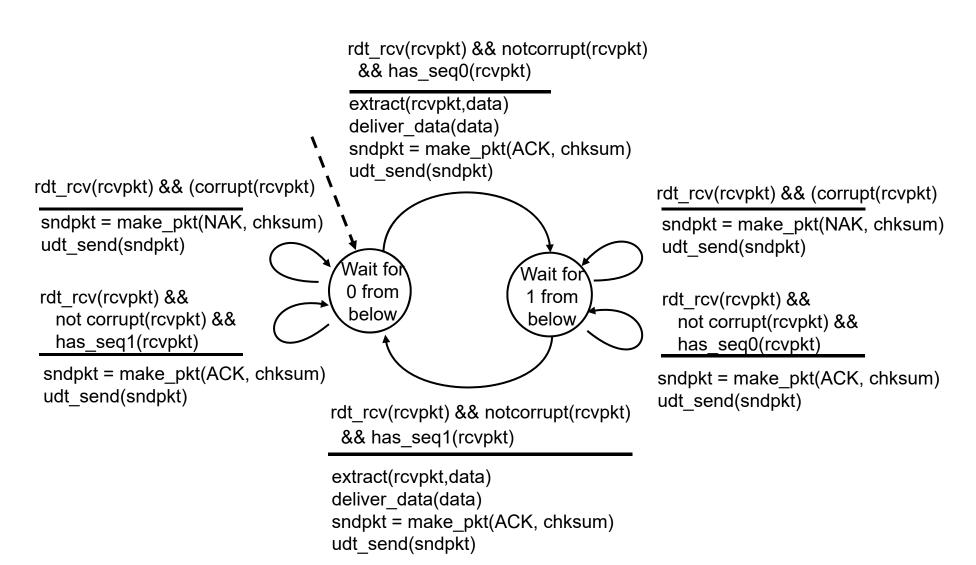
stop and wait

sender sends one packet, then waits for receiver response

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



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



#### rdt2.1: discussion

#### sender:

- seq # added to pkt
- two seq. #'s (0,1) will suffice. Why?
- must check if received ACK/NAK corrupted

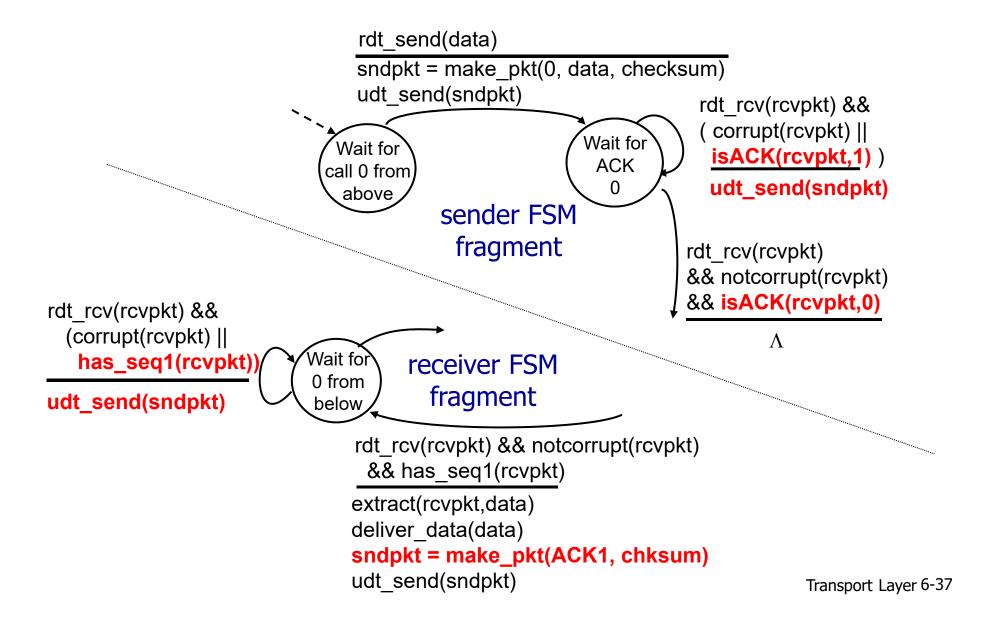
#### receiver:

- must check if received packet is duplicate
  - state indicates whether0 or I is expected pktseq #
- note: receiver can not know if its last ACK/NAK received OK at sender

## rdt2.2: a NAK-free protocol

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

### rdt2.2: sender, receiver fragments



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

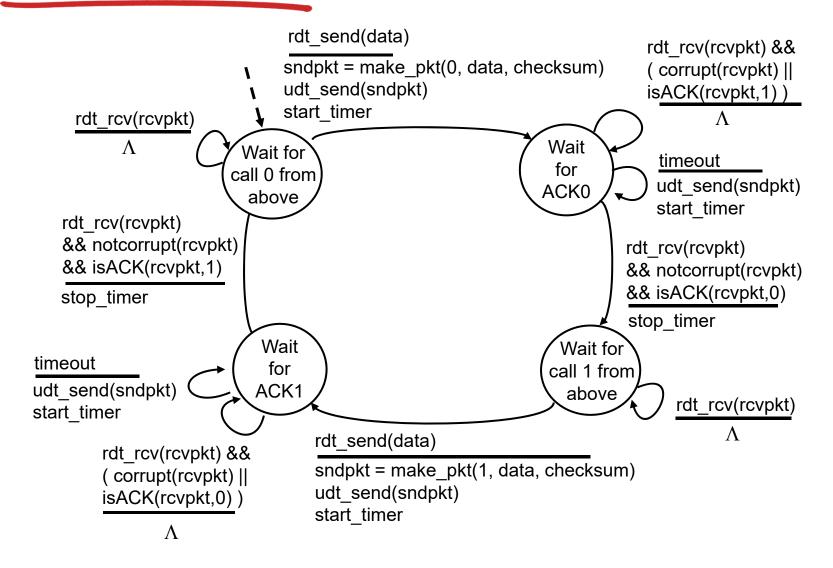
#### new assumption:

underlying channel can also lose packets (data, ACKs)

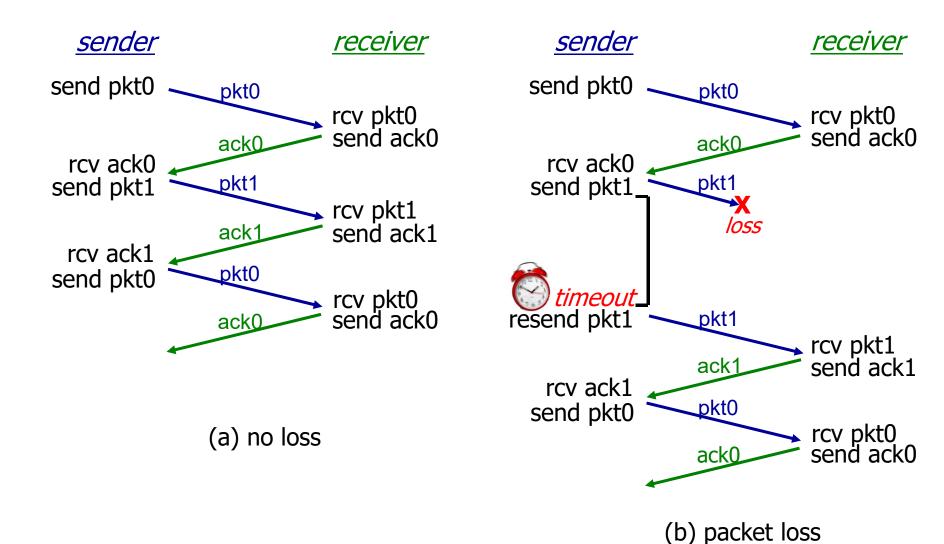
checksum, seq. #,
 ACKs, retransmissions
 will be of help ... but
 not enough

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

#### rdt3.0 sender

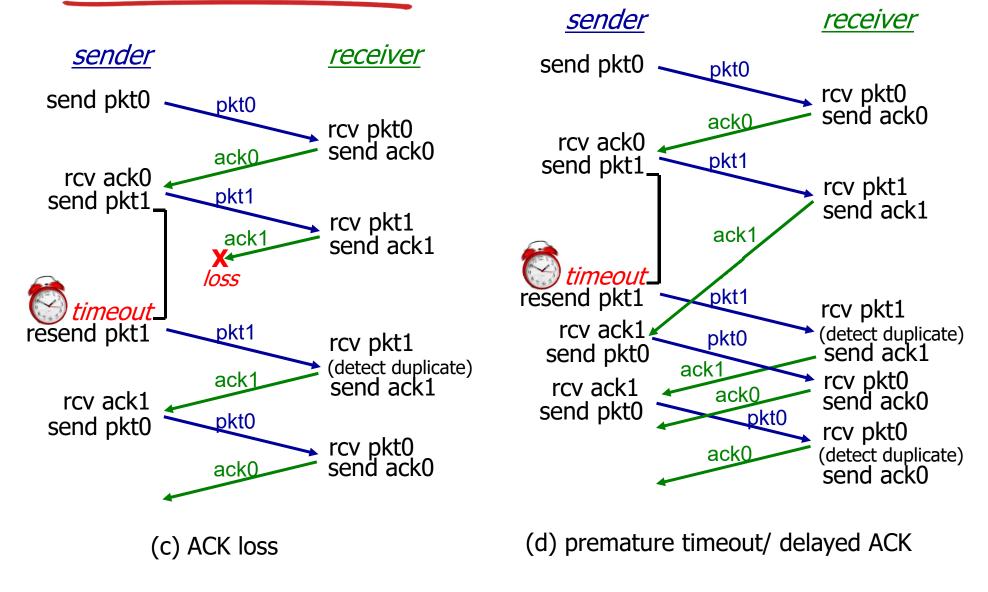


# rdt3.0 in action



Transport Layer 6-40

### rdt3.0 in action



# <u>Outline</u>

- 6.1 transport-layer services
- 6.2 multiplexing and demultiplexing
- 6.3 connectionless transport: UDP
- 6.4 principles of reliable data transfer

- 6.5 connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- 6.6 principles of congestion control
- 6.7 TCP congestion control

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

- point-to-point:
  - one sender, one receiver
- reliable, in-order byte steam:
  - no "message boundaries"
- pipelined:
  - TCP congestion and flow control set window size
- socket door TCP send buffer segment segment application reads data socket door

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

## TCP segment structure

32 bits URG: urgent data counting dest port # source port # (generally not used) by bytes sequence number of data ACK: ACK # (not segments!) acknowledgement number valid head not len used UAP receive window PSH: push data now # bytes (generally not used) cheeksum Urg data pointer rcvr willing to accept RST, SYN, FIN: options (variable length) connection estab (setup, teardown commands) application data Internet (variable length) checksum<sup>2</sup> (as in UDP)

# TCP seq. numbers, ACKs

#### sequence numbers:

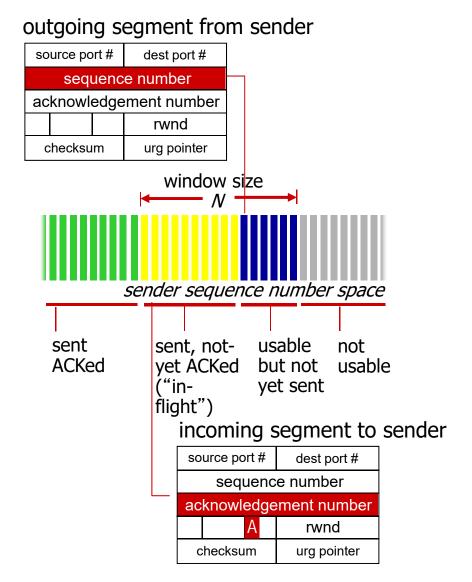
• byte stream "number" of first byte in segment's data

#### acknowledgements:

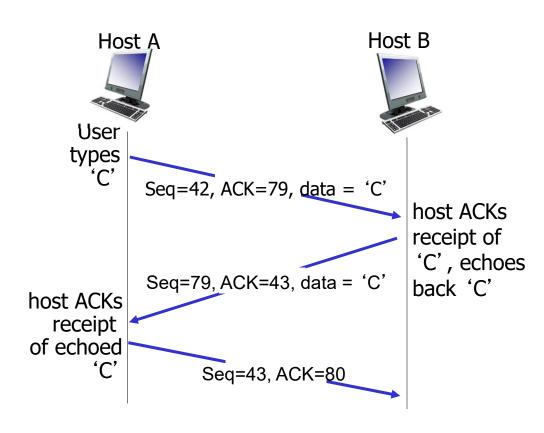
- seq # of next byte expected from other side
- cumulative ACK

Q: how receiver handles out-of-order segments

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



# TCP seq. numbers, ACKs



simple telnet scenario

# TCP round trip time, timeout

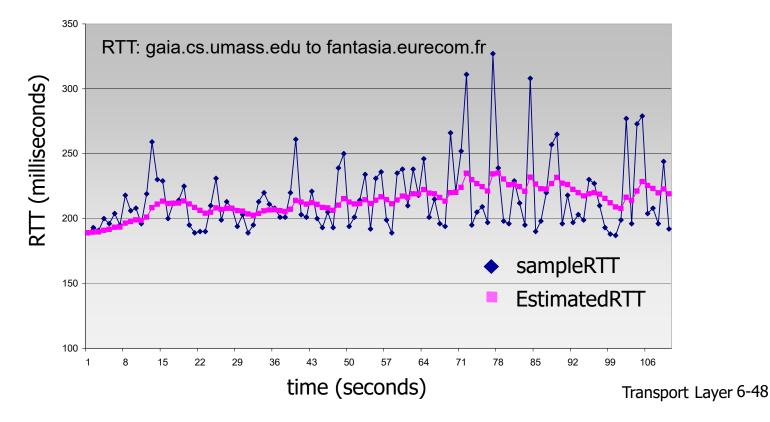
- Q: how to set TCP timeout value?
- longer than RTT
  - but RTT varies
- too short: premature timeout, unnecessary retransmissions
- too long: slow reaction to segment loss

- Q: how to estimate RTT?
- SampleRTT: measured time from segment transmission until ACK receipt
  - ignore retransmissions
- SampleRTT will vary, want estimated RTT "smoother"
  - average several recent measurements, not just current SampleRTT

# TCP round trip time, timeout

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

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



# TCP round trip time, timeout

- \* timeout interval: EstimatedRTT plus "safety margin"
  - large variation in EstimatedRTT -> larger safety margin
- estimate SampleRTT deviation from EstimatedRTT:

```
DevRTT = (1-\beta)*DevRTT +

\beta*|SampleRTT-EstimatedRTT|

(typically, \beta = 0.25)
```

TimeoutInterval = EstimatedRTT + 4\*DevRTT



estimated RTT "safety margin"

# **Outline**

- 6.1 transport-layer services
- 6.2 multiplexing and demultiplexing
- 6.3 connectionless transport: UDP
- 6.4 principles of reliable data transfer

- 6.5 connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- 6.6 principles of congestion control
- 6.7 TCP congestion control

## TCP reliable data transfer

- TCP creates rdt service on top of IP's unreliable service
  - pipelined segments
  - cumulative acks
  - single retransmission timer
- retransmissions triggered by:
  - timeout events
  - duplicate acks

let's initially consider simplified TCP sender:

- ignore duplicate acks
- ignore flow control, congestion control

### TCP sender events:

#### data rcvd from app:

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

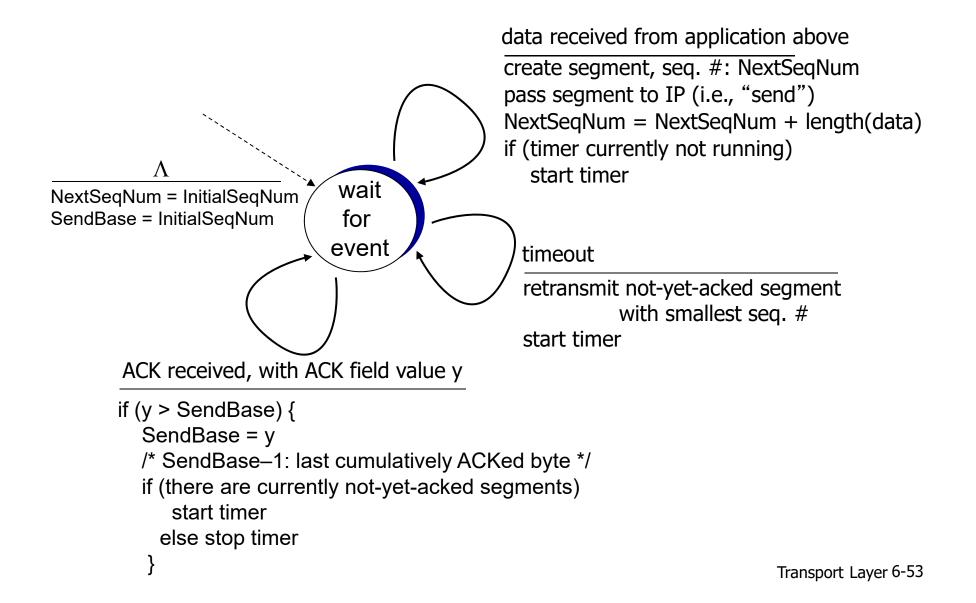
#### timeout:

- retransmit segment that caused timeout
- restart timer

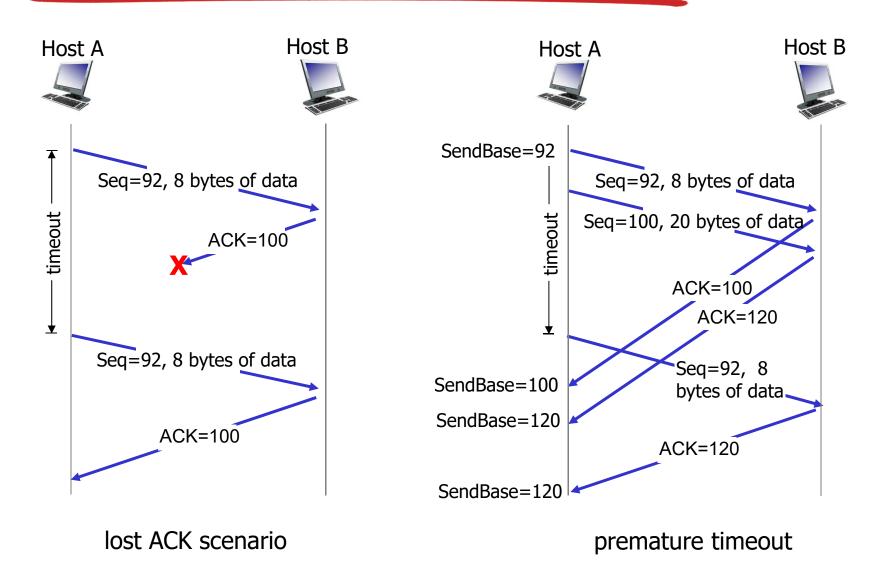
#### ack rcvd:

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

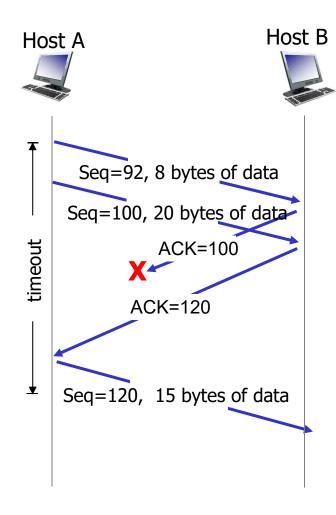
# TCP sender (simplified)



### TCP: retransmission scenarios



### TCP: retransmission scenarios



cumulative ACK

# TCP ACK generation [RFC 1122, RFC 2581]

event at receiver	TCP receiver action
arrival of in-order segment with expected seq #. All data up to expected seq # already ACKed	delayed ACK. Wait up to 500ms for next segment. If no next segment, send ACK
arrival of in-order segment with expected seq #. Sender has 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 <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

## TCP fast retransmit

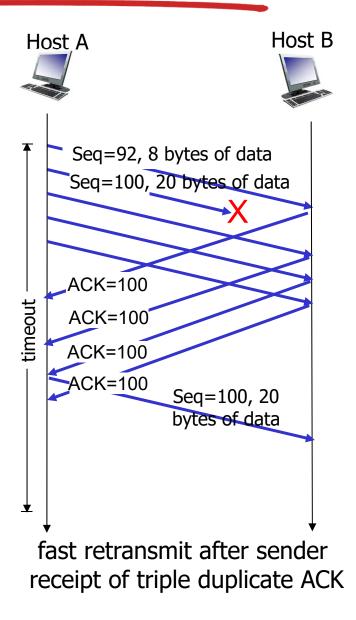
- time-out period often relatively long:
  - long delay before resending lost packet
- detect lost segments via duplicate ACKs.
  - sender often sends many segments backto-back
  - if segment is lost, there will likely be many duplicate ACKs.

#### TCP fast retransmit

if sender receives 3
ACKs for same data
("triple duplicate ACKs"),
resend unacked
segment with smallest
seq #

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

## TCP fast retransmit



## Outline

- 6.1 transport-layer services
- 6.2 multiplexing and demultiplexing
- 6.3 connectionless transport: UDP
- 6.4 principles of reliable data transfer

- 6.5 connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- 6.6 principles of congestion control
- 6.7 TCP congestion control

## TCP flow control

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

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

### application process application OS TCP socket receiver buffers **TCP** code ΙP code from sender

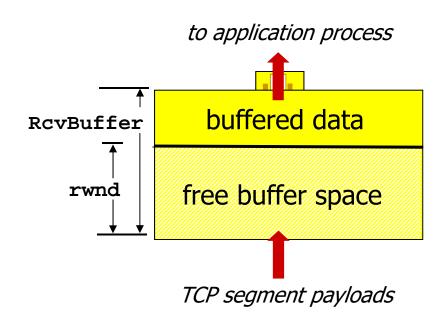
receiver protocol stack

#### flow control

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

## 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



receiver-side buffering

# **Outline**

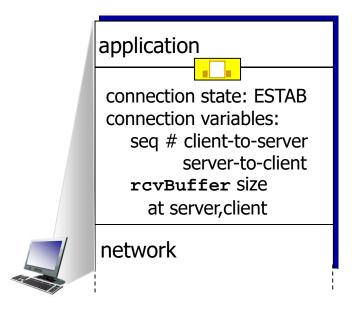
- 6.1 transport-layer services
- 6.2 multiplexing and demultiplexing
- 6.3 connectionless transport: UDP
- 6.4 principles of reliable data transfer

- 6.5 connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- 6.6 principles of congestion control
- 6.7 TCP congestion control

### Connection Management

#### before exchanging data, sender/receiver "handshake":

- agree to establish connection (each knowing the other willing to establish connection)
- agree on connection parameters



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

```
application

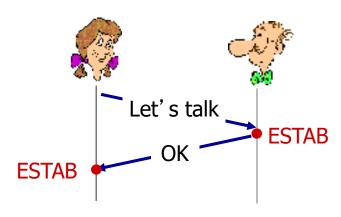
connection state: ESTAB
connection Variables:
  seq # client-to-server
        server-to-client
   rcvBuffer size
   at server,client

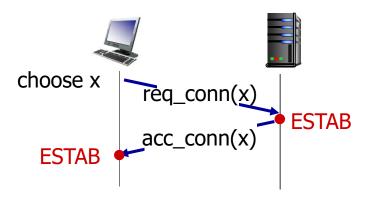
network
```

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

### Agreeing to establish a connection

#### 2-way handshake:

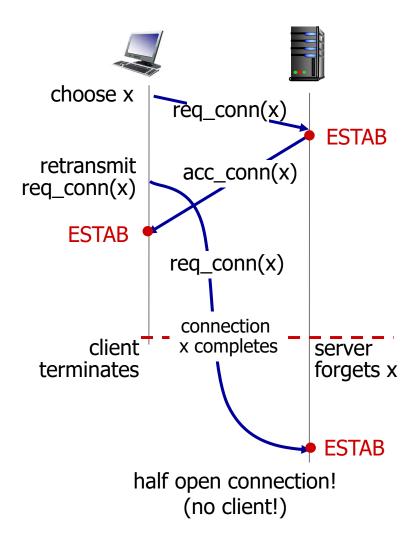




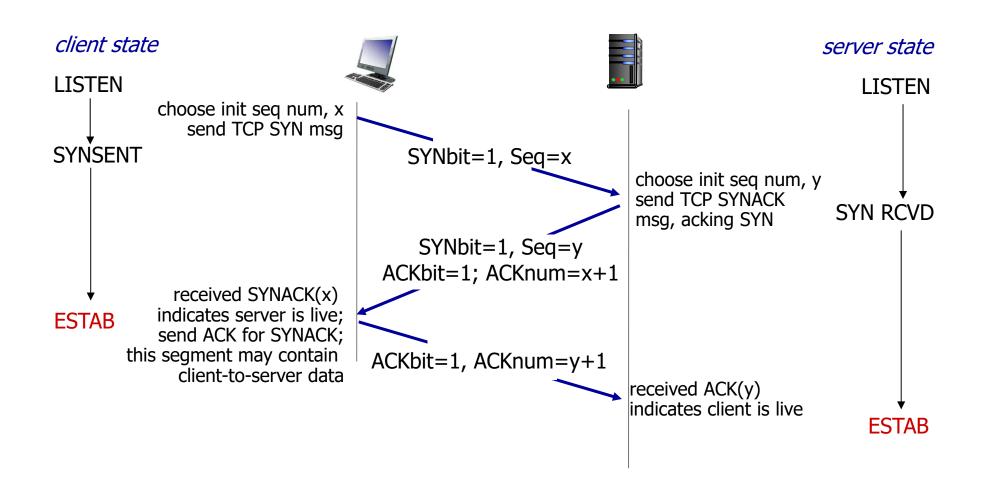
- Q: will 2-way handshake always work in network?
- variable delays
- retransmitted messages (e.g. req\_conn(x)) due to message loss
- message reordering
- can't "see" other side

### Agreeing to establish a connection

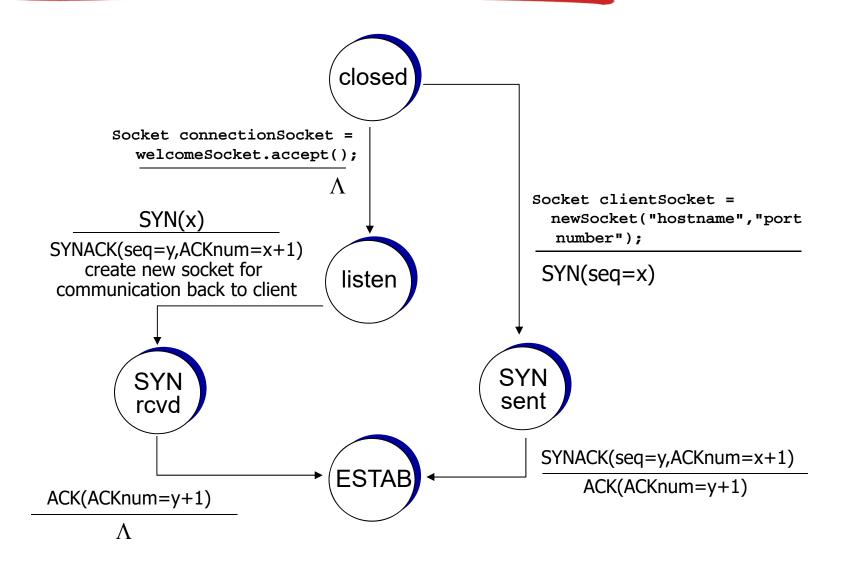
#### 2-way handshake failure scenario:



### TCP 3-way handshake



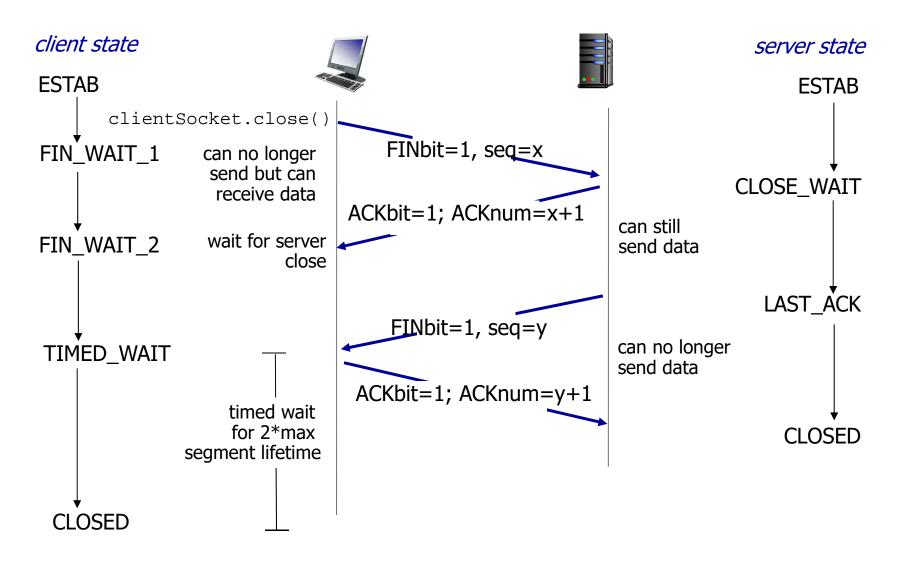
## TCP 3-way handshake: FSM



# TCP: closing a connection

- client, server each close their side of connection
  - send TCP segment with FIN bit = I
- 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



## Outline

- 6.1 transport-layer services
- 6.2 multiplexing and demultiplexing
- 6.3 connectionless transport: UDP
- 6.4 principles of reliable data transfer

- 6.5 connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- 6.6 principles of congestion control
- 6.7 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!

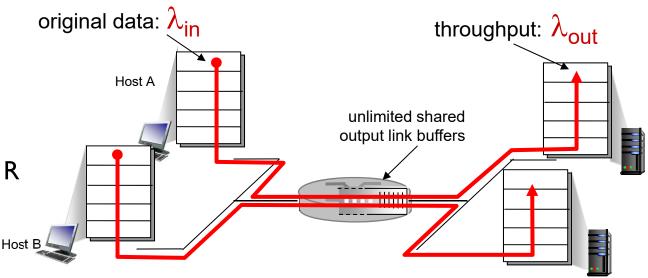
## Causes/costs of congestion: scenario I

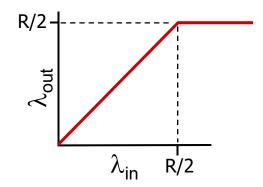
two senders, two receivers

one router, infinite buffers

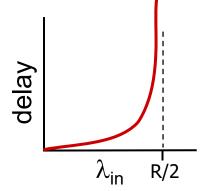
output link capacity: R

no retransmission



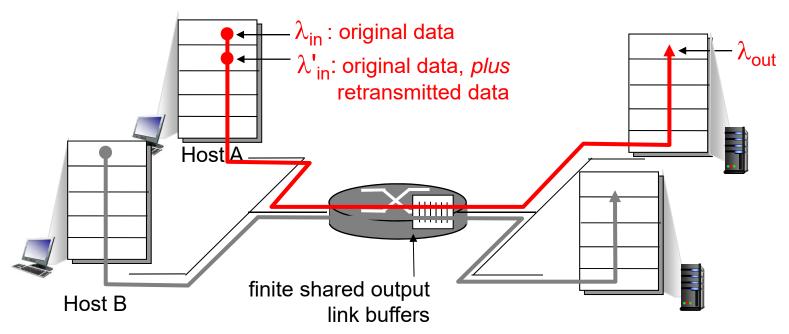


maximum per-connection throughput: R/2



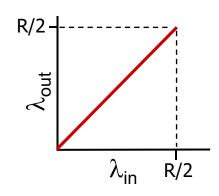
• large delays as arrival rate,  $\lambda_{in}$ , approaches capacity

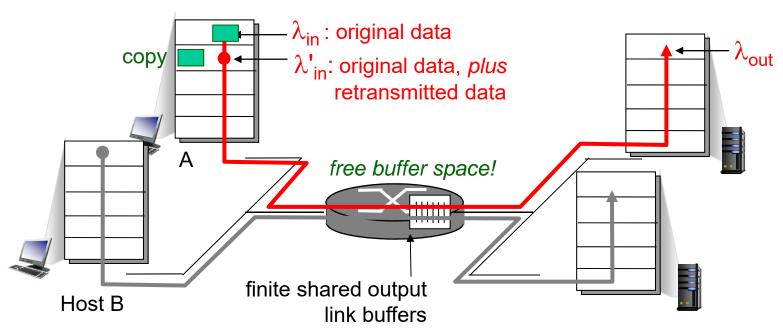
- one router, finite buffers
- sender retransmission of timed-out packet
  - application-layer input = application-layer output:  $\lambda_{\text{in}}$  =  $\lambda_{\text{out}}$
  - transport-layer input includes retransmissions :  $\lambda_{in} \geq \lambda_{in}$



## idealization: perfect knowledge

 sender sends only when router buffers available

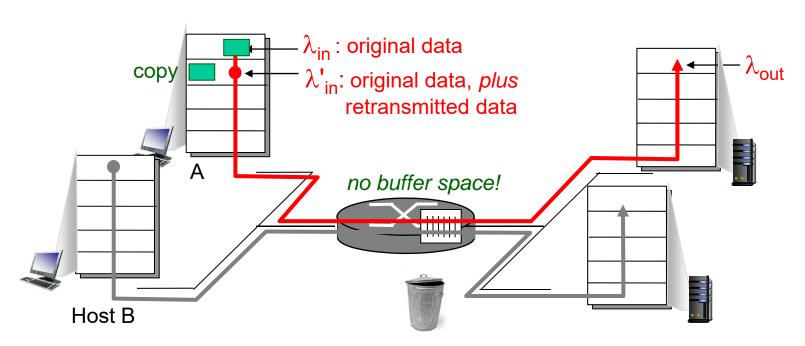




#### Idealization: known loss

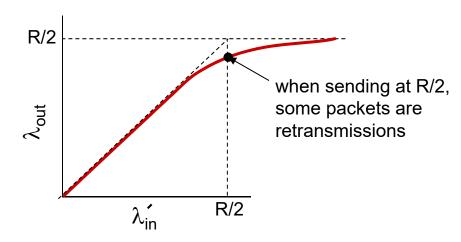
packets can be lost, dropped at router due to full buffers

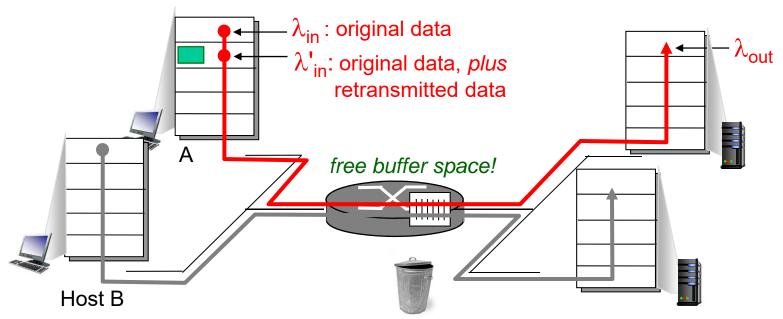
sender only resends if packet known to be lost



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

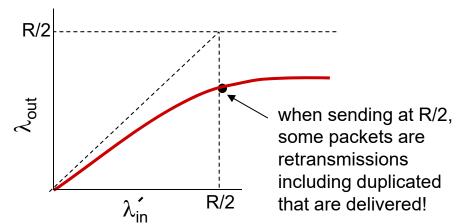
sender only resends if packet known to be lost

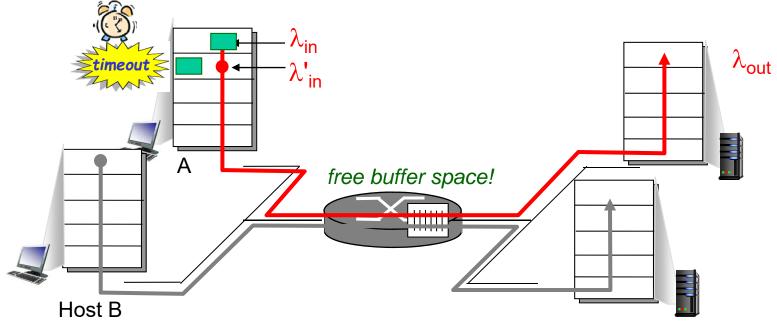




### Realistic: duplicates

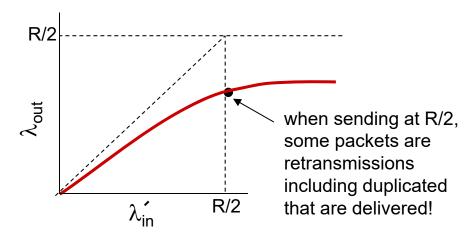
- packets can be lost, dropped at router due to full buffers
- sender times out prematurely, sending two copies, both of which are delivered





### 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

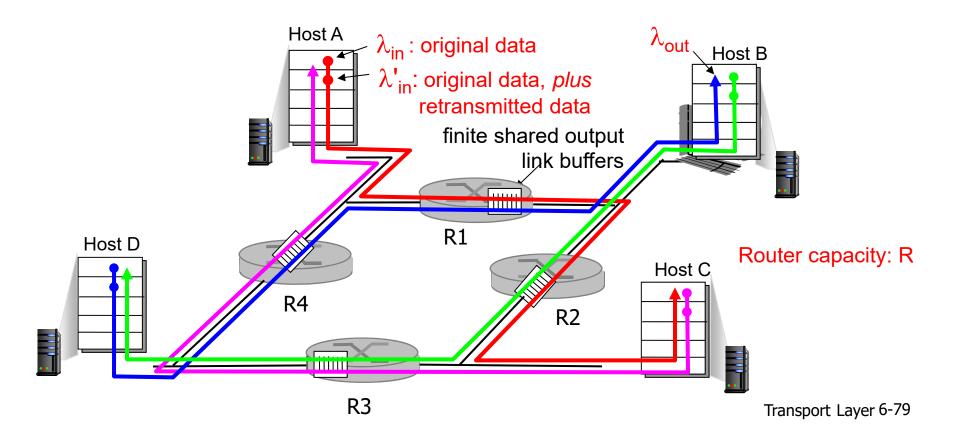


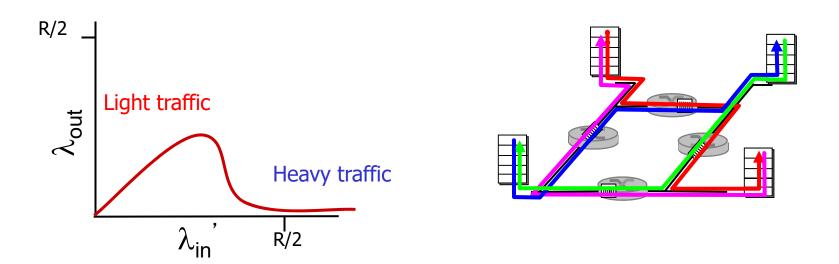
### "costs" of congestion:

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

- four senders
- multihop paths
- timeout/retransmit

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





### another "cost" of congestion:

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

### Approaches towards congestion control

### two broad approaches towards congestion control:

## end-end congestion control:

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

## network-assisted congestion control:

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

### Outline

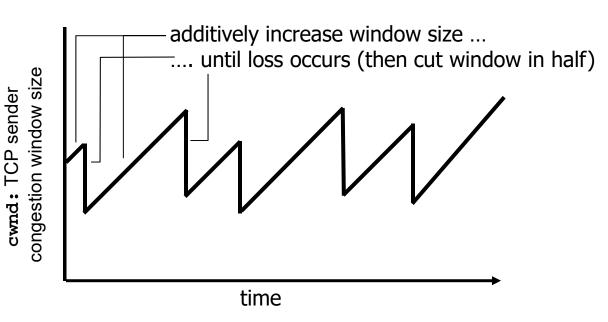
- 6.1 transport-layer services
- 6.2 multiplexing and demultiplexing
- 6.3 connectionless transport: UDP
- 6.4 principles of reliable data transfer

- 6.5 connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- 6.6 principles of congestion control
- 6.7 TCP congestion control

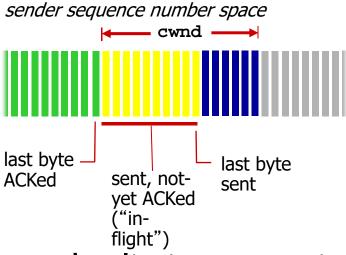
## TCP congestion control: additive increase multiplicative decrease

- \* approach: sender increases transmission rate (window size), probing for usable bandwidth, until loss occurs
  - additive increase: increase cwnd by I MSS every RTT until loss detected
  - multiplicative decrease: cut cwnd in half after loss

AIMD saw tooth behavior: probing for bandwidth



## TCP Congestion Control: details



sender limits transmission:

$$\begin{array}{ccc} \texttt{LastByteSent-} & \leq & \texttt{cwnd} \\ \texttt{LastByteAcked} & \end{array}$$

 cwnd is dynamic, function of perceived network congestion

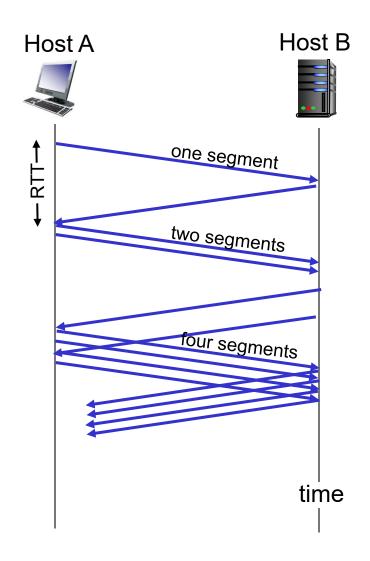
### TCP sending rate:

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

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

### TCP Slow Start

- when connection begins, increase rate exponentially until first loss event:
  - initially cwnd = I MSS
  - double cwnd every RTT
  - done by incrementing cwnd for every ACK received
- summary: initial rate is slow but ramps up exponentially fast



## TCP: detecting, reacting to loss

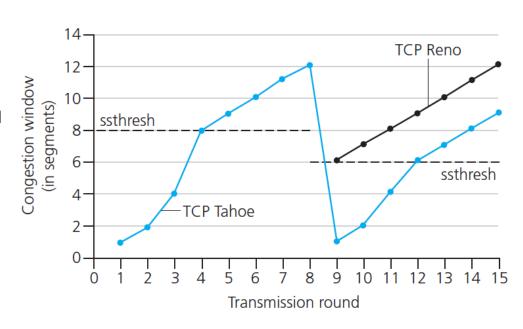
- loss indicated by timeout:
  - cwnd set to I MSS;
  - window then grows exponentially (as in slow start) to threshold, then grows linearly
- loss indicated by 3 duplicate ACKs: TCP RENO
  - dup ACKs indicate network not capable of delivering some segments
  - cwnd is cut in half window then grows linearly
- TCP Tahoe always sets cwnd to I (timeout or 3 duplicate acks)

### TCP: switching from slow start to CA

Q: when should the exponential increase switch to linear?

A: when cwnd gets to 1/2 of its value when congestion was last encountered.

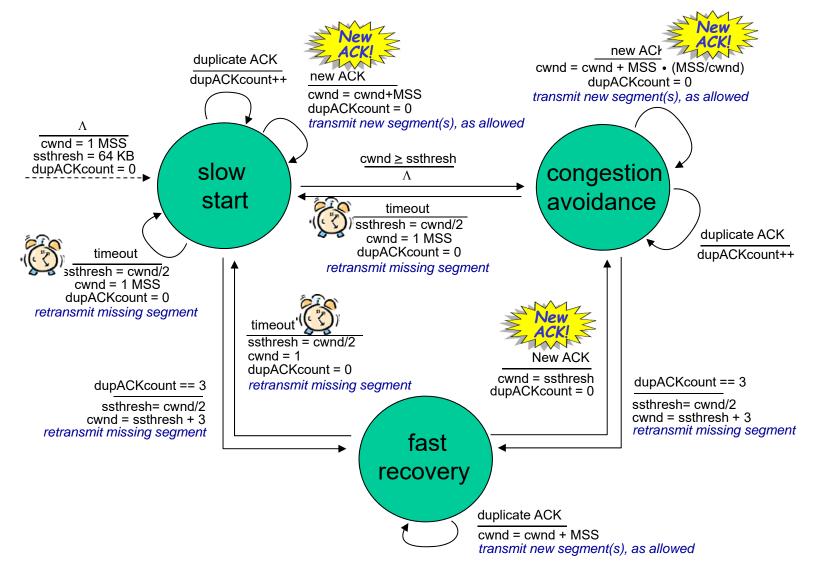
#### **Congestion Avoidance**



### **Implementation:**

- variable ssthresh
- on loss event, ssthresh
  is set to 1/2 of cwnd
  when the loss event
  occurred

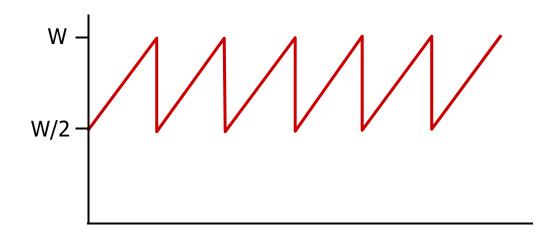
### Summary: TCP Congestion Control



## TCP throughput

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

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



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

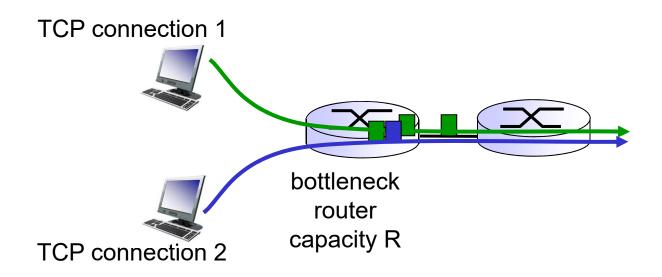
- example: I500 byte segments, I00ms RTT, want
   I0 Gbps throughput
- requires W = 83,333 in-flight segments
- throughput in terms of segment loss probability, L [Mathis 1997]:

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

- → to achieve 10 Gbps throughput, need a loss rate of L =  $2 \cdot 10^{-10} a$  very small loss rate!
- new versions of TCP for high-speed

### TCP Fairness

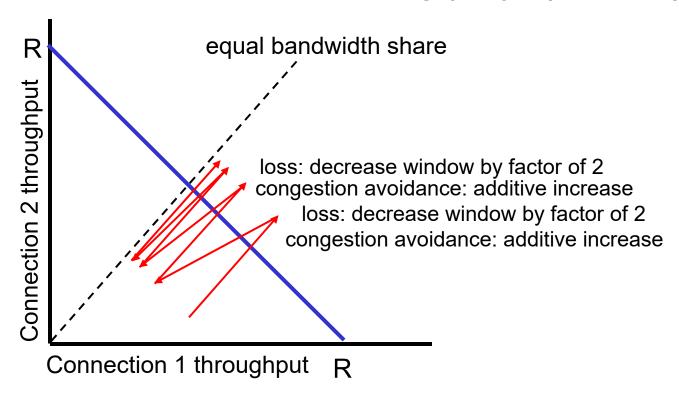
fairness goal: if K TCP sessions share same bottleneck link of bandwidth R, each should have average rate of R/K



## Why is TCP fair?

#### two competing sessions:

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



## 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 I TCP, gets rate R/10
  - new app asks for 11 TCPs, gets R/2

## T6: summary

- principles behind transport layer services:
  - multiplexing, demultiplexing
  - reliable data transfer
  - flow control
  - congestion control
- instantiation, implementation in the Internet
  - UDP
  - TCP