# Master Course
# Computer Networks
# IN2097

**Prof. Dr.-Ing. Georg Carle**
**Christian Grothoff, Ph.D.**
*Lecturer today: Dr. Nils Kammenhuber*

**Chair for Network Architectures and Services**

**Institut für Informatik**
**Technische Universität München**
**http://www.net.in.tum.de**

Technische Universität München

# Chapter 3: Transport Layer

Our goals:

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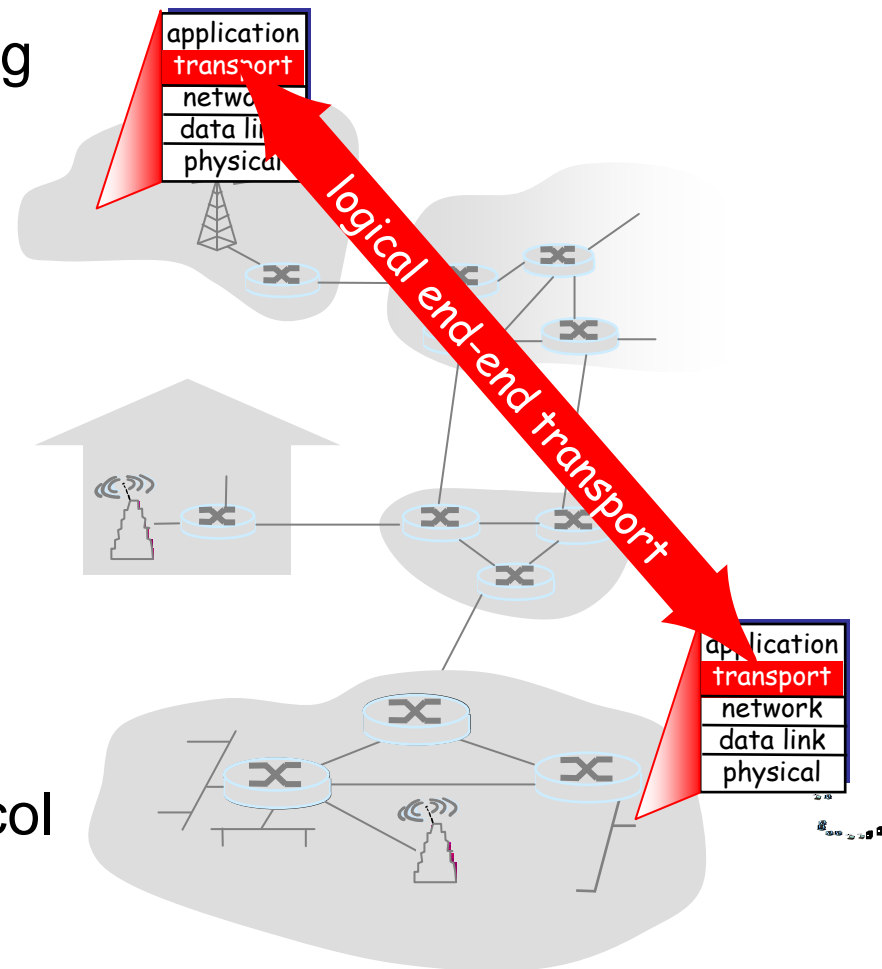  - TCP: connection-oriented transport
  - TCP congestion control

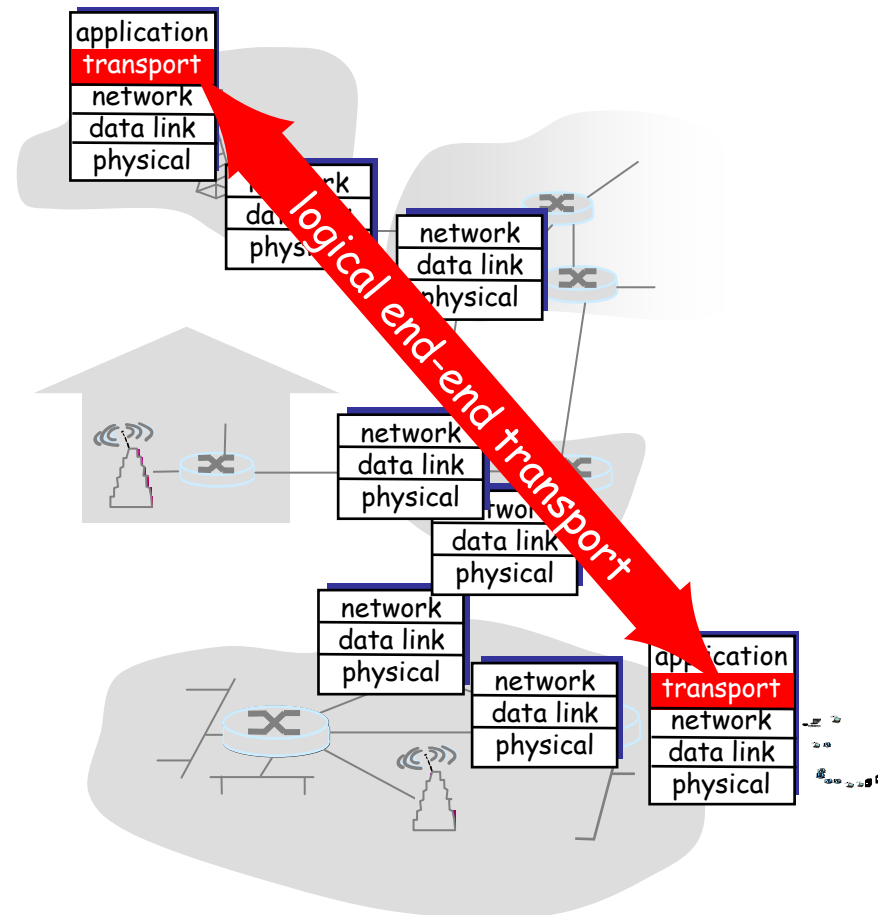# Chapter 3 outline

# 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

# 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
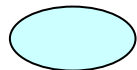
<u>Demultiplexing at rcv host:</u>

delivering received segments
to correct socket

<u>Multiplexing at send host:</u>

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

▭ = socket      ⬭ = process

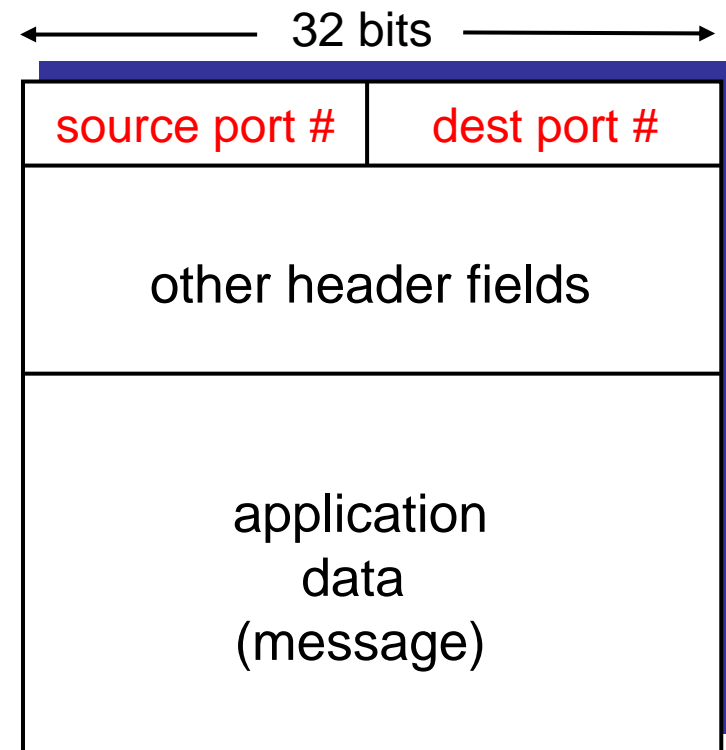| application | P3 | | P1 | application | P2 | | P4 | application |
| transport | | | | transport | | | | transport |
| network | | | | network | | | | network |
| link | | | | link | | | | link |
| physical | | | | physical | | | | physical |

host 1              host 2              host 3

# How demultiplexing works

❑ **host receives IP datagrams**

- each datagram has source IP address, destination IP address

- each datagram carries 1 transport-layer segment

- each segment has source, destination port number

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

32 bits

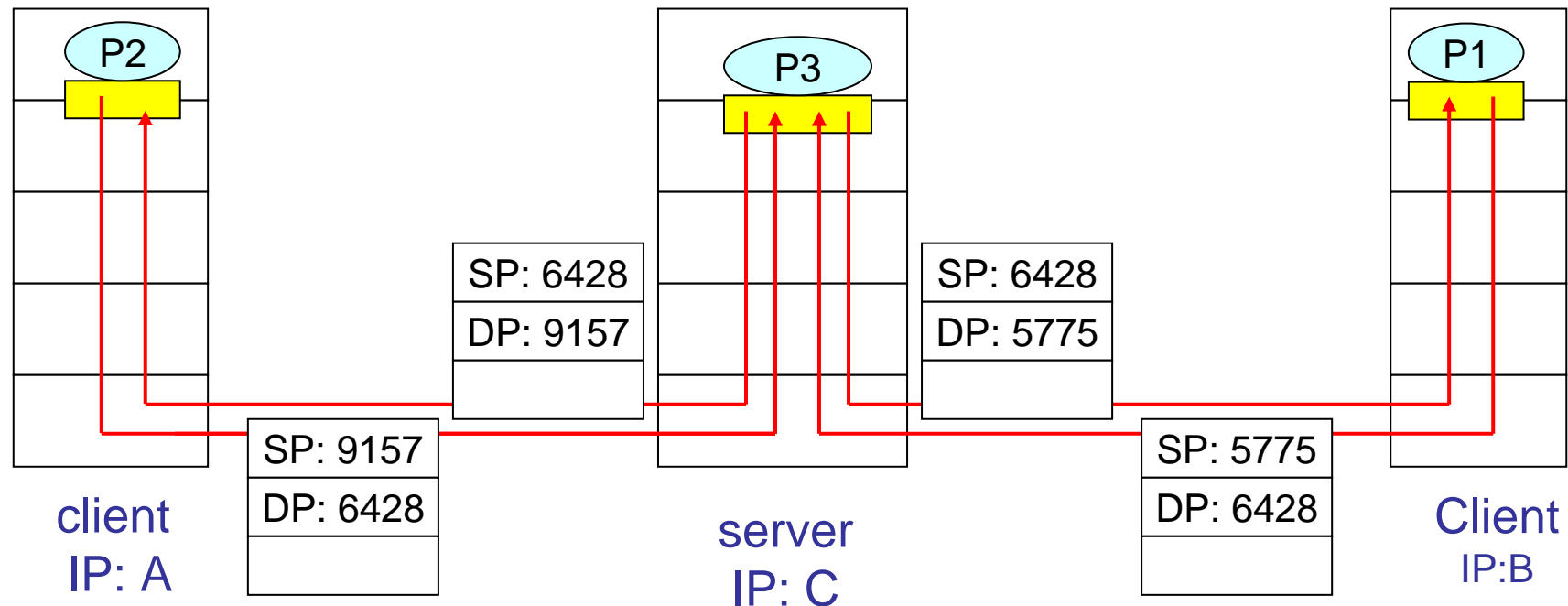| source port # | dest port # |
|---|---|
| other header fields | |
| application data (message) | |

TCP/UDP segment format

# Connectionless demultiplexing

❑ Create sockets with port numbers:

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

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

❑ 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

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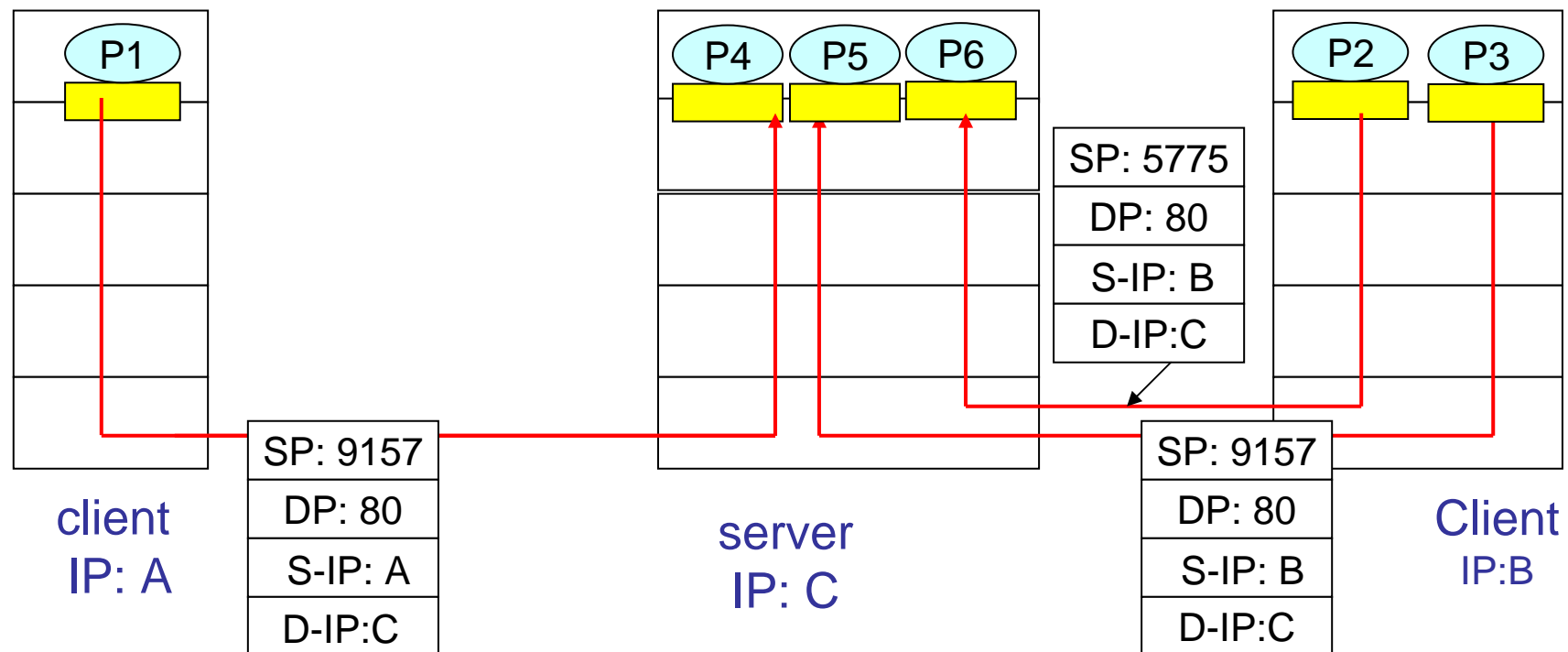


Source Port (SP) provides "return address"

# Connection-oriented demux

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

P1

P4  P5  P6

P2  P3

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

client
IP: A

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

server
IP: C

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

Client
IP:B

P1

P4

P2   P3

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

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

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

client
IP: A

server
IP: C

Client
IP:B

# Chapter 3 outline

- 3.1 Transport-layer services

- 3.2 Multiplexing and demultiplexing

- **3.3 Connectionless transport: UDP**

- -

- 3.5 Connection-oriented transport: TCP

  - segment structure

  - reliable data transfer

  - flow control

  - connection management

- -

- 3.7 TCP congestion control

# UDP: User Datagram Protocol [RFC 768]

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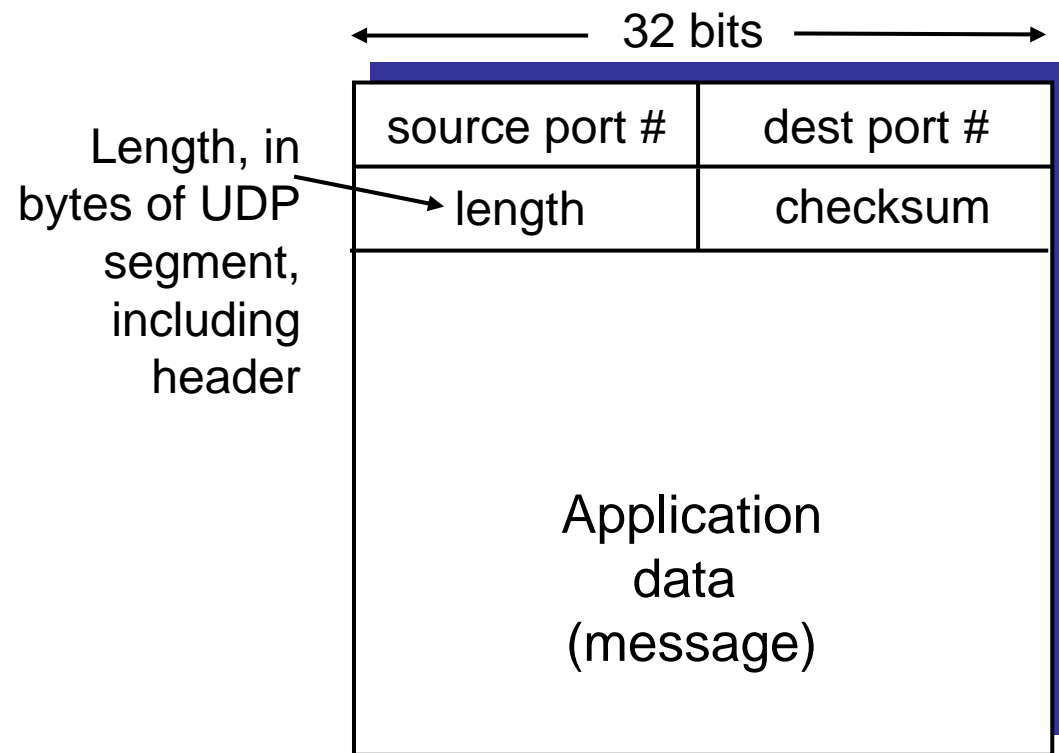  - ▪ each UDP segment handled independently of others

Why is there a UDP?

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

# UDP: more

- often used for streaming multimedia apps
  - loss tolerant
  - rate sensitive
- other UDP uses
  - DNS
  - SNMP
- reliable transfer over UDP: add reliability at application layer
  - application-specific error recovery!

← 32 bits →

| source port # | dest port # |
|---|---|
| length | checksum |

Length, in bytes of UDP segment, including header → length

Application
data
(message)

UDP segment format

# UDP checksum

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

Sender:

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

Receiver:

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

- ❑ Note
  - ▪ When adding numbers, a carryout from the most significant bit needs to be added to the result
- ❑ Example: add two 16-bit integers

```
              1 1 1 0 0 1 1 0 0 1 1 0 0 1 1 0
              1 1 0 1 0 1 0 1 0 1 0 1 0 1 0 1
              ─────────────────────────────────
wraparound  ⓵ 1 0 1 1 1 0 1 1 1 0 1 1 1 0 1 1
              ─────────────────────────────────
       sum    1 0 1 1 1 0 1 1 1 0 1 1 1 1 0 0
  checksum    0 1 0 0 0 1 0 0 0 1 0 0 0 0 1 1
```

# Chapter 3 outline
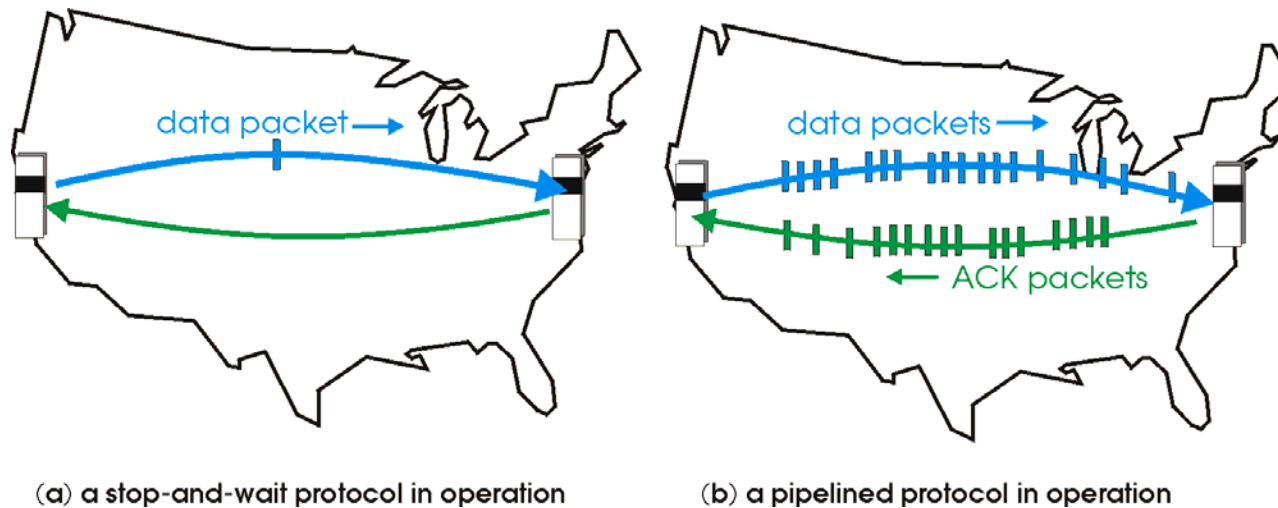
- 3.1 Transport-layer services
- 3.2 Multiplexing and demultiplexing
- 3.3 Connectionless transport: UDP
- **3.4 Principles of reliable data transfer**
- 3.5 Connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- 3.6 Principles of congestion control
- 3.7 TCP congestion control

# Pipelined protocols

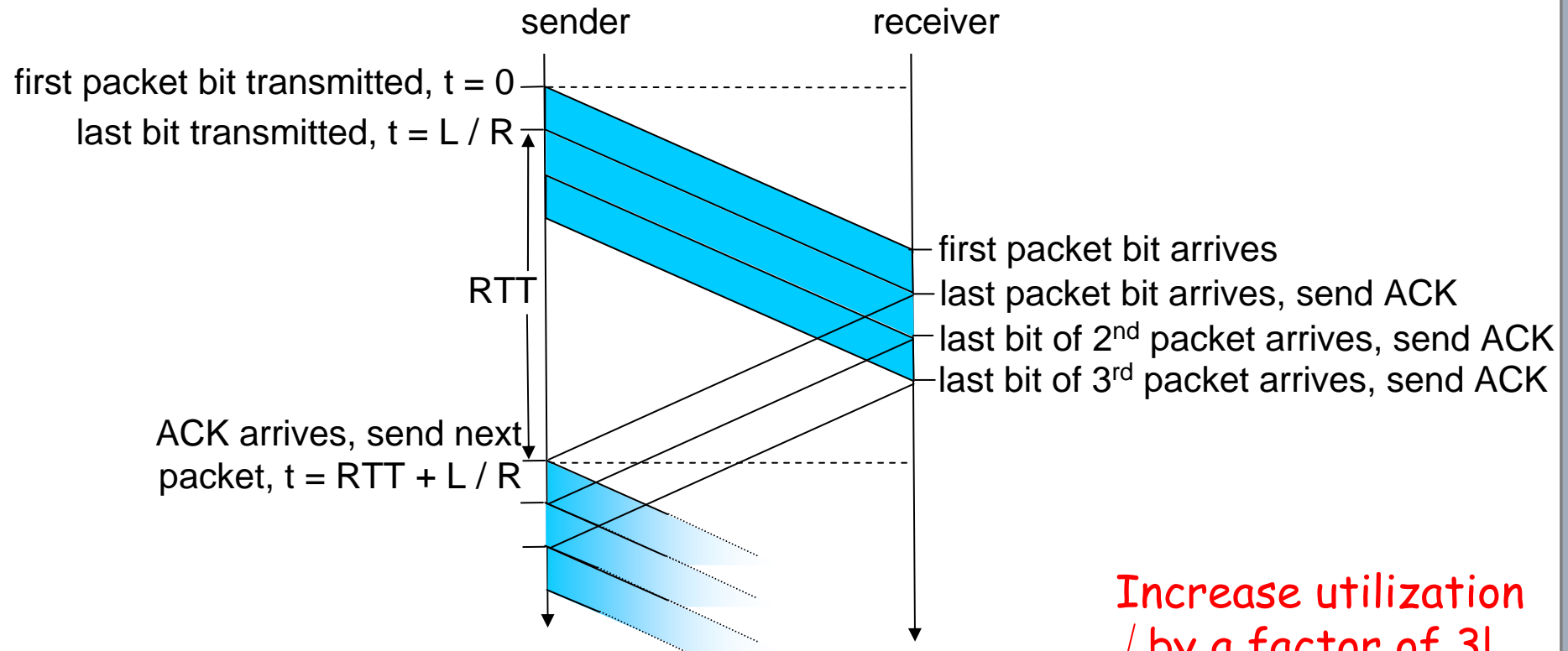Pipelining: sender allows multiple, "in-flight", yet-to-be-acknowledged pkts

- range of sequence numbers must be increased
- buffering at sender and/or receiver



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

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

# Pipelining: increased utilization

sender       receiver

first packet bit transmitted, t = 0

last bit transmitted, t = L / R

RTT

first packet bit arrives

last packet bit arrives, send ACK

last bit of 2$^{nd}$ packet arrives, send ACK

last bit of 3$^{rd}$ packet arrives, send ACK

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

Increase utilization
by a factor of 3!

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

# Go-Back-N

Sender:

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



- ❑ ACK(n): ACKs all pkts up to, including seq # n - "cumulative ACK"
  - ■ may receive duplicate ACKs (see receiver)
- ❑ timer for each in-flight pkt
- ❑ *timeout(n):* retransmit pkt n and all higher seq # pkts in window
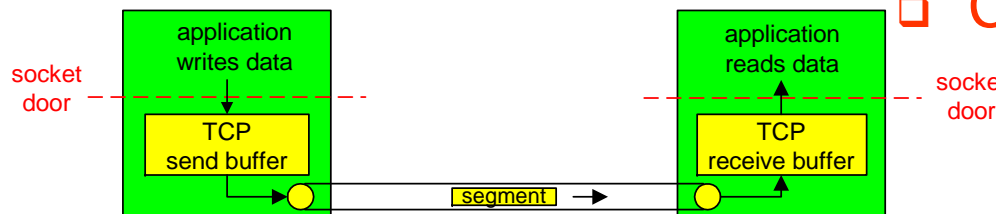
# Chapter 3 outline

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

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

socket door

application writes data

TCP send buffer

application reads data

TCP receive buffer

socket door

segment

# TCP segment structure

32 bits

URG: urgent data
(generally not used)

ACK: ACK #
valid

PSH: push data now
(generally not used)

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

Internet
checksum
(as in UDP)

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

| head len | not used | U | A | P | R | S | F | Receive window |
|---|---|---|---|---|---|---|---|---|
| checksum | | | | | | | | Urg data pnter |

Options (variable length)

application
data
(variable length)

counting
by bytes
of data
(not segments!)

# bytes
rcvr willing
to accept

# TCP seq. #'s and ACKs
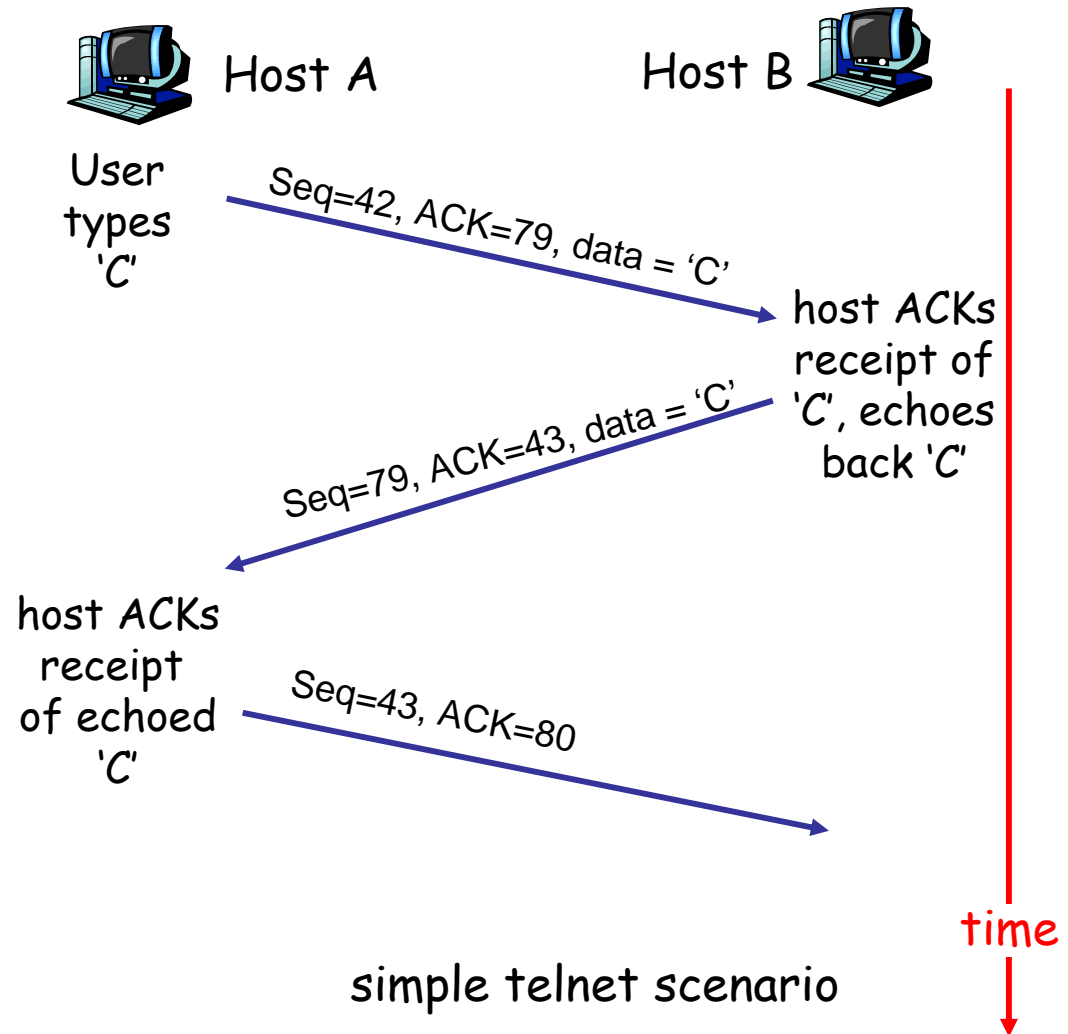
<u>Seq. #'s:</u>

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

<u>ACKs:</u>

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

Q: how receiver handles out-of-order segments

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

Host A                                      Host B

User
types
'C'
                    Seq=42, ACK=79, data = 'C'
                                                host ACKs
                                                receipt of
                                                'C', echoes
                    Seq=79, ACK=43, data = 'C'  back 'C'

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

time

simple telnet scenario

# TCP Round Trip Time and Timeout

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

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

**Q:** how to estimate RTT?

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

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

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

# Example RTT estimation:

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

## Setting the timeout

- [ ] **`EstimtedRTT`** plus "safety margin"
  - large variation in **`EstimatedRTT -> `** larger safety margin
- [ ] first estimate of how much SampleRTT deviates from EstimatedRTT:

$$\texttt{DevRTT = (1-}\beta\texttt{)*DevRTT +}$$
$$\beta\texttt{*|SampleRTT-EstimatedRTT|}$$

$$\texttt{(typically, } \beta \texttt{ = 0.25)}$$

Then set timeout interval:

`TimeoutInterval = EstimatedRTT + 4*DevRTT`

# Chapter 3 outline

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

# TCP reliable data transfer

- TCP creates rdt service on top of IP's unreliable service
- Pipelined segments
- Cumulative acks
- TCP uses single retransmission timer

- Retransmissions are triggered by:
  - timeout events
  - duplicate acks
- Initially consider simplified TCP sender:
  - ignore duplicate acks
  - ignore flow control, congestion control

# TCP sender events:

**data rcvd from app:**

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

**timeout:**

- retransmit segment that caused timeout
- restart timer

**Ack rcvd:**

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

# TCP sender (simplified)

```
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 */
```

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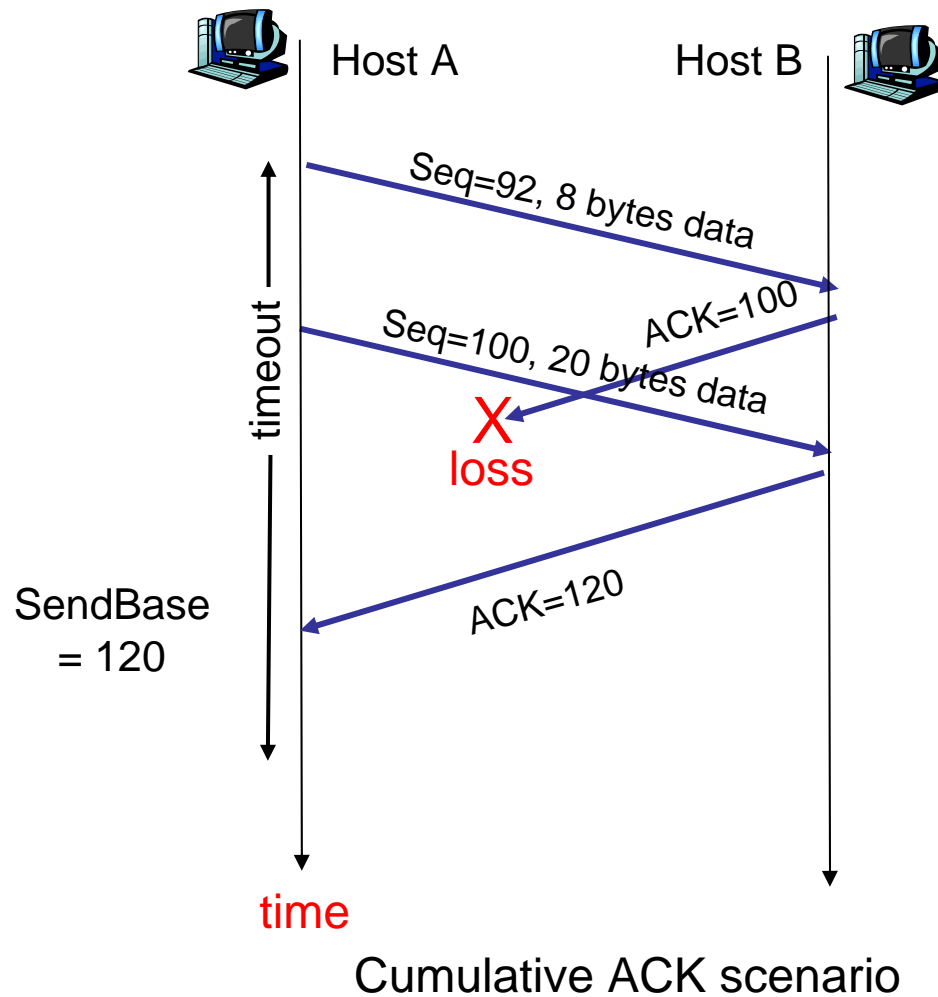
# TCP: retransmission scenarios



Host A                          Host B

Seq=92, 8 bytes data

ACK=100

timeout

X
loss

Seq=92, 8 bytes data

ACK=100

SendBase
= 100

time

lost ACK scenario

---

Host A                          Host B

Seq=92, 8 bytes data

Seq=100, 20 bytes data

ACK=100

ACK=120

Seq=92 timeout

Seq=92, 8 bytes data

Sendbase
= 100
SendBase
= 120

Seq=92 timeout

ACK=120

SendBase
= 120

time

premature timeout

Host A          Host B

Seq=92, 8 bytes data

ACK=100

Seq=100, 20 bytes data

X
loss

timeout

SendBase
= 120

ACK=120

time

Cumulative ACK scenario

# TCP ACK generation [RFC 1122, RFC 2581]

| Event at Receiver | TCP Receiver action |
|---|---|
| Arrival of in-order segment with expected seq #. All data up to expected seq # already ACKed | Delayed ACK. Wait up to 500ms for next segment. If no next segment, send ACK |
| Arrival of in-order segment with expected seq #. One other segment has ACK pending | Immediately send single cumulative ACK, ACKing both in-order segments |
| Arrival of out-of-order segment higher-than-expect seq. # . Gap detected | Immediately send *duplicate ACK*, indicating seq. # of next expected byte |
| Arrival of segment that partially or completely fills gap | Immediate send ACK, provided that segment starts at lower end of gap |

# Fast  Retransmit

- Time-out period  often relatively long:
    - long delay before resending lost packet
- Detect lost segments via duplicate ACKs.
    - Sender often sends many segments back-to-back
    - If segment is lost, there will likely be many duplicate ACKs.

- If sender receives 3 ACKs for the same data, it supposes that segment after ACKed data was lost:
    - <u>fast retransmit:</u> resend segment before timer expires

# 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

# Chapter 3 outline

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

# TCP Flow Control

❑ receive side of TCP connection
has a receive buffer:



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

❑ app process may be slow at reading
from buffer

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

# TCP Flow control: how it works



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

❑ spare room in buffer

= `RcvWindow`

= `RcvBuffer-[LastByteRcvd -
  LastByteRead]`

❑ Rcvr advertises spare room by including value of `RcvWindow` in segments

❑ Sender limits unACKed data to `RcvWindow`

- guarantees receive buffer doesn't overflow

# Chapter 3 outline

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

# TCP Connection Management

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

- initialize TCP variables:
    - seq. #s
    - buffers, flow control info (e.g. `RcvWindow`)
- *client:* connection initiator

    `Socket clientSocket = new`

    `Socket("hostname","port number");`
- *server:* contacted by client

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

## Three way handshake:

<u>Step 1:</u> client host sends TCP SYN segment to server

- specifies initial seq #
- no data

<u>Step 2:</u> server host receives SYN, replies with SYNACK segment

- server allocates buffers
- specifies server initial seq. #

<u>Step 3:</u> client receives SYNACK, replies with ACK segment, which may contain data

# TCP Connection Management (cont.)

## Closing a connection:

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

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

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

Step 3: client receives FIN, replies with ACK.

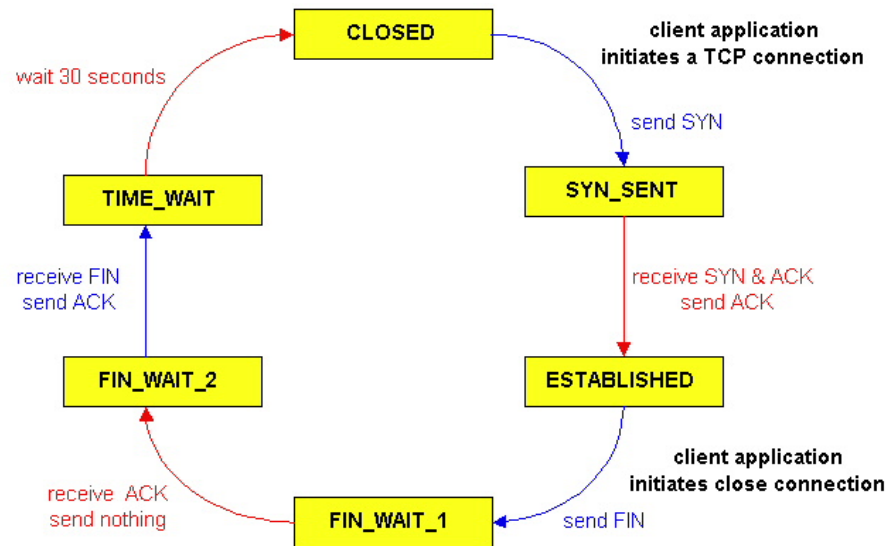- Enters "timed wait" - will respond with ACK to received FINs

Step 4: server, receives ACK. Connection closed.
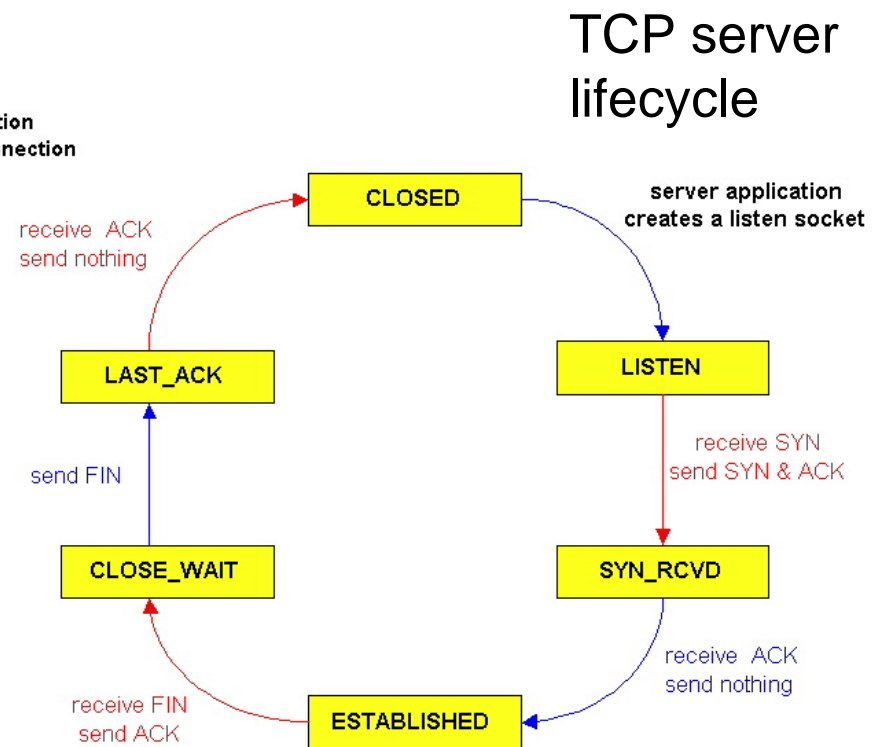
Note: with small modification, can handle simultaneous FINs.

TCP client lifecycle

TCP server lifecycle
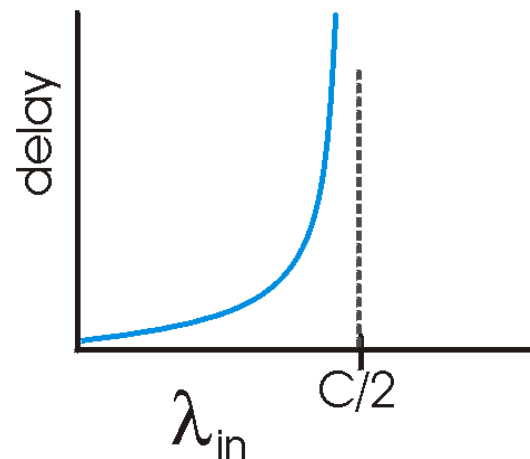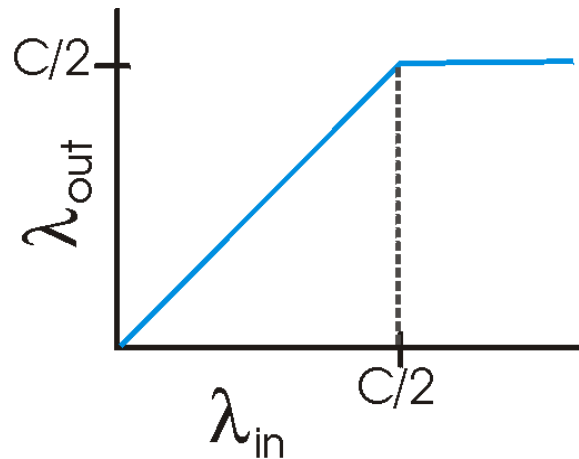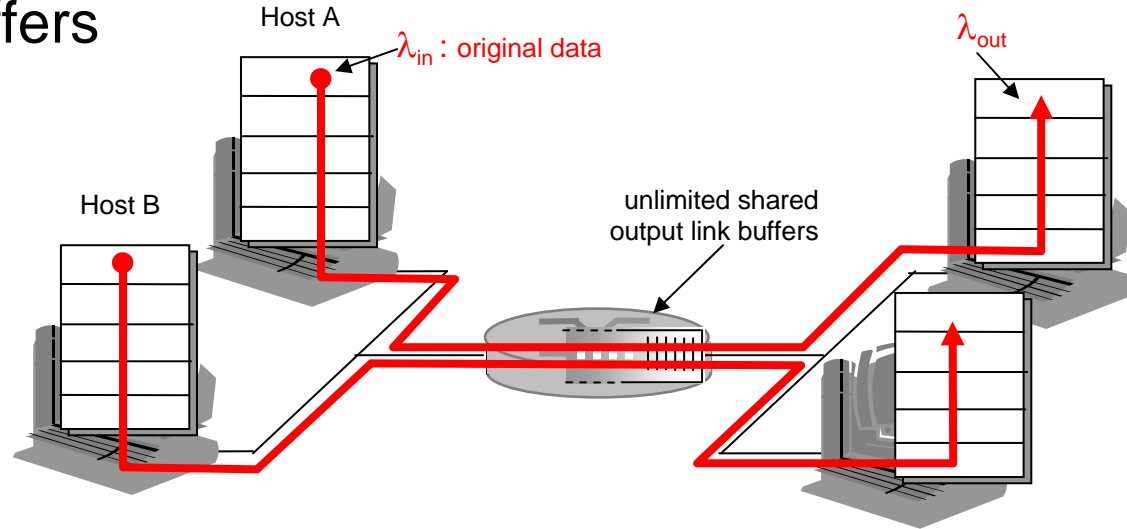
# Chapter 3 outline

# 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!

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



Host A

$\lambda_{in}$ : original data

$\lambda_{out}$

Host B

unlimited shared output link buffers



- large delays when congested
- maximum achievable throughput

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

Host A

$\lambda_{in}$ : original data

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

$\lambda_{out}$

Host B

finite shared output link buffers
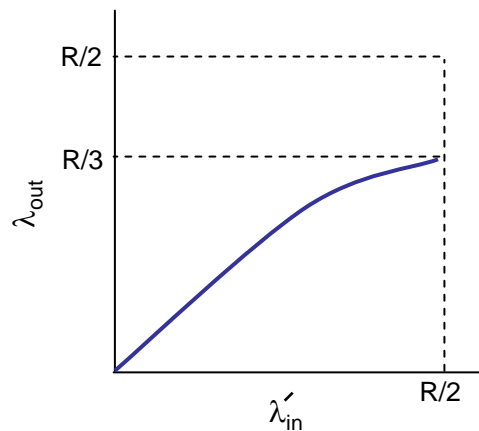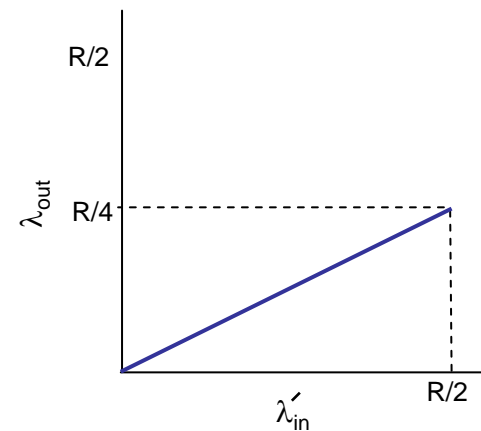
- always: $\lambda_{in} = \lambda_{out}$ (goodput)

- "perfect" retransmission only when loss: $\lambda'_{in} > \lambda_{out}$

- retransmission of delayed (not lost) packet makes $\lambda'_{in}$ larger (than perfect case) for same $\lambda_{out}$



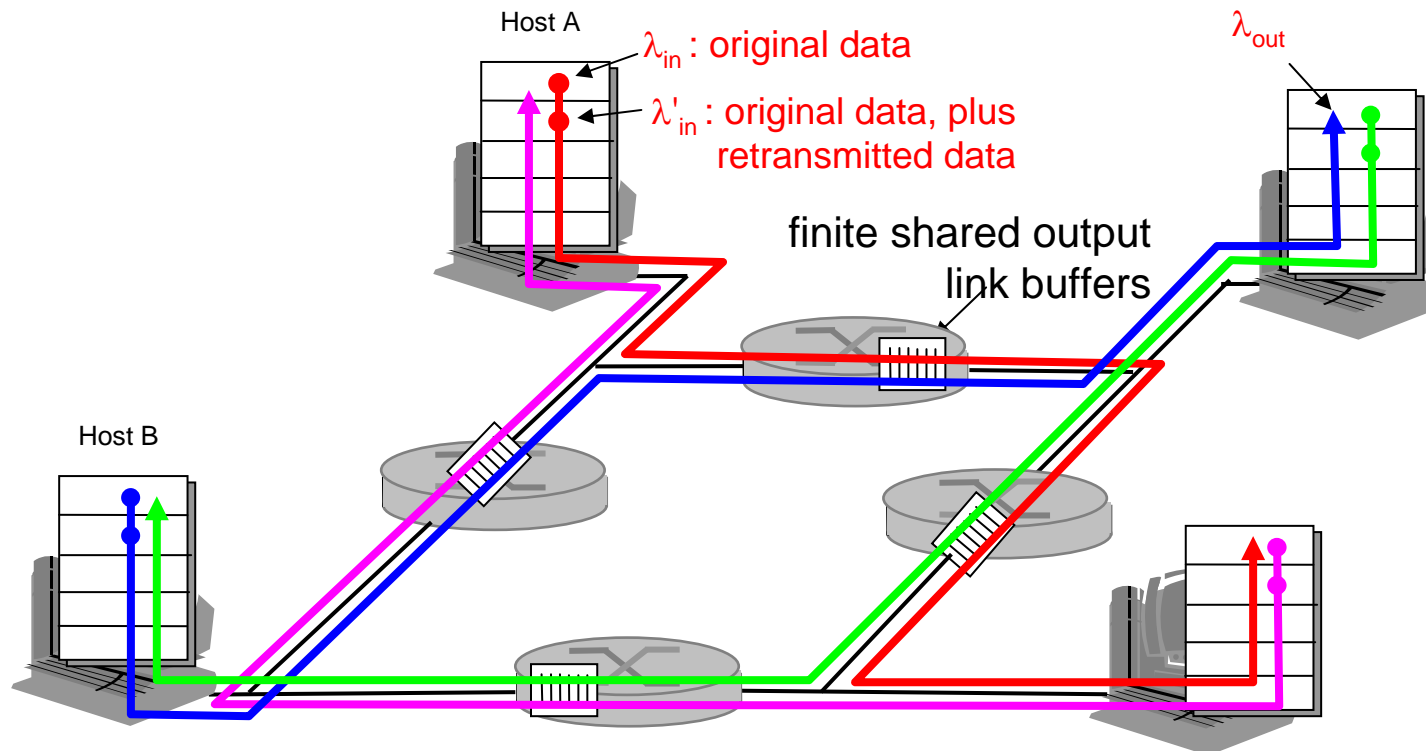a.          b.          c.

"costs" of congestion:

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

- four senders
- multihop paths
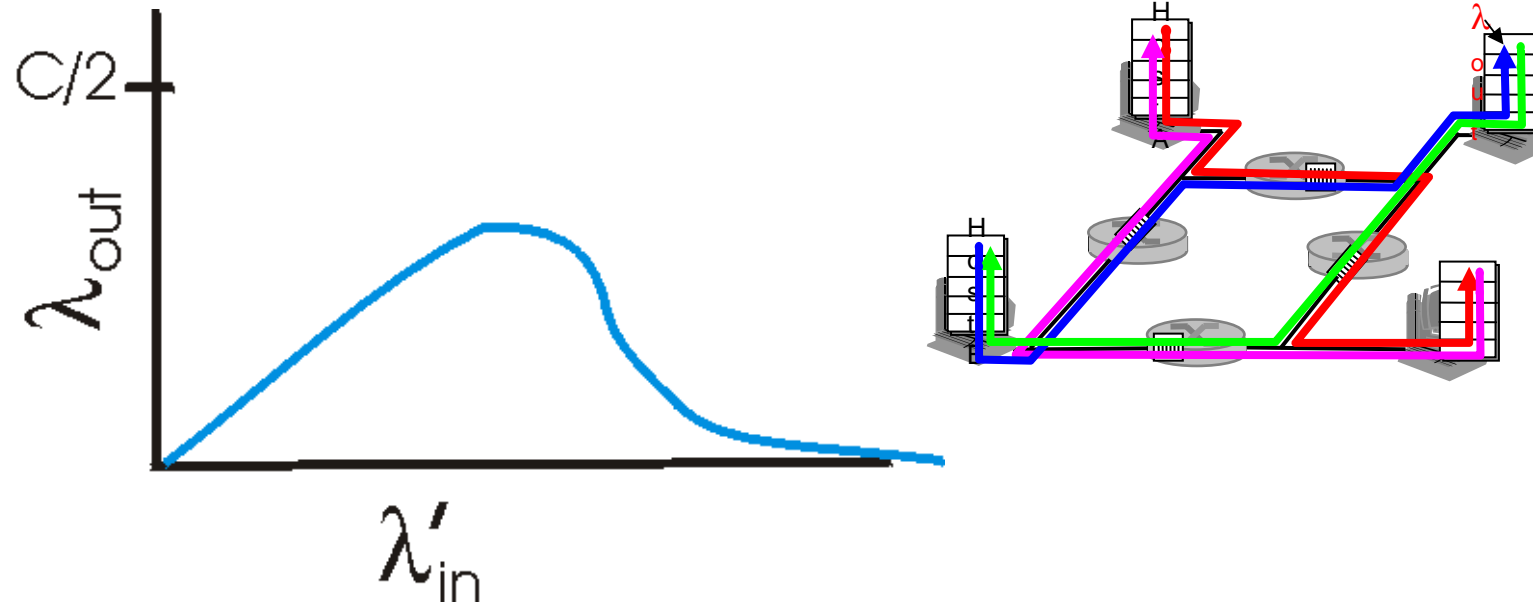- timeout/retransmit

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

Host A

$\lambda_{in}$ : original data

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

$\lambda_{out}$

finite shared output link buffers

Host B

# Causes/costs of congestion: scenario 3



Another "cost" of congestion:

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

# Approaches towards congestion control

Two broad approaches towards congestion control:

End-end congestion control:

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

Network-assisted congestion control:

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

# Case study: ATM ABR congestion control

ABR: available bit rate:

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

RM (resource management) cells:

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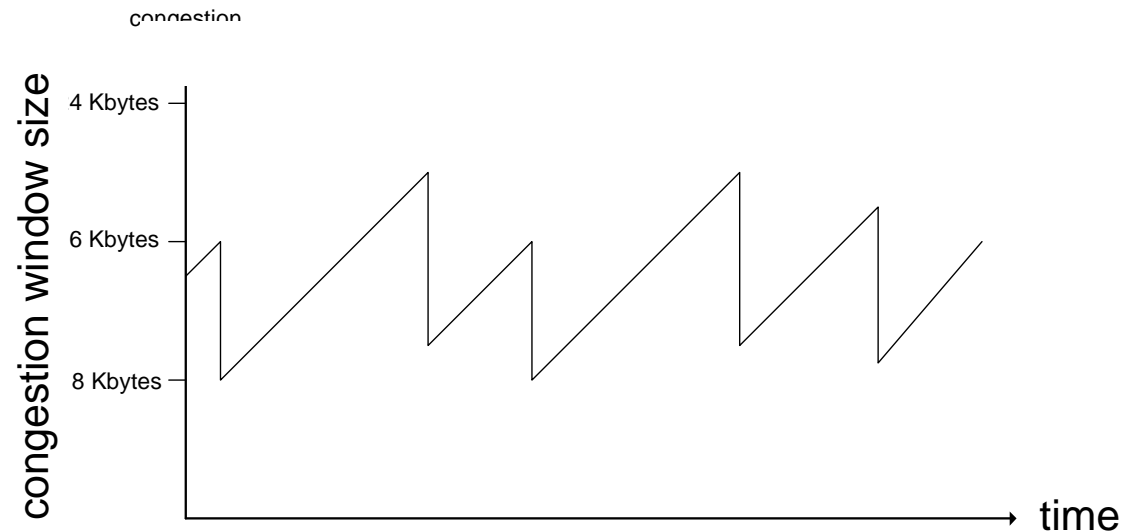- RM cells returned to sender by receiver, with bits intact

# Chapter 3 outline

❑ *Approach:* increase transmission rate (window size), probing for usable bandwidth, until loss occurs

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

   ▪ *multiplicative decrease*: cut **CongWin** in half after loss

Saw tooth behavior: probing for bandwidth

congestion

congestion window size

4 Kbytes —

6 Kbytes —

8 Kbytes —

time

□ sender limits transmission:

**LastByteSent-LastByteAcked**

$$\leq \texttt{CongWin}$$

□ Roughly,

$$\text{rate} = \frac{\texttt{CongWin}}{\text{RTT}} \ \text{Bytes/sec}$$

□ **CongWin** is dynamic, function of perceived network congestion

## How does sender perceive congestion?

□ loss event = timeout *or* 3 duplicate acks

□ TCP sender reduces rate (**CongWin**) after loss event

## three mechanisms:
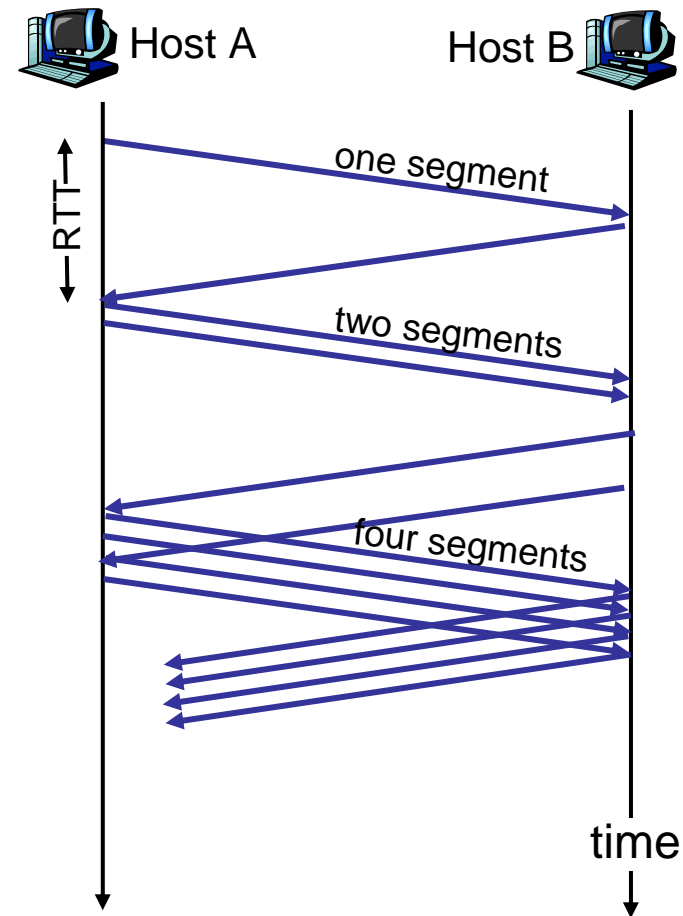
- AIMD
- slow start
- conservative after timeout events

❑ When connection begins,
  `CongWin` = 1 MSS

  ▪ Example: MSS = 500
    bytes & RTT = 200 msec

  ▪ initial rate = 20 kbps

❑ available bandwidth may be
  >> MSS/RTT

  ▪ desirable to quickly ramp
    up to respectable rate

❑ When connection begins,
  increase rate exponentially
  fast until first loss event

# TCP Slow Start (more)

❑ When connection begins, increase rate exponentially until first loss event:
   ▪ double `CongWin` every RTT
   ▪ done by incrementing `CongWin` 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

❑ After 3 dup ACKs:

- ▪ `CongWin` is cut in half

- ▪ window then grows linearly

❑ <u>But</u> after timeout event:

- ▪ `CongWin` instead set to 1 MSS;

- ▪ window then grows exponentially

- ▪ to a threshold, then grows linearly

Philosophy:

❑ 3 dup ACKs indicates network capable of delivering some segments

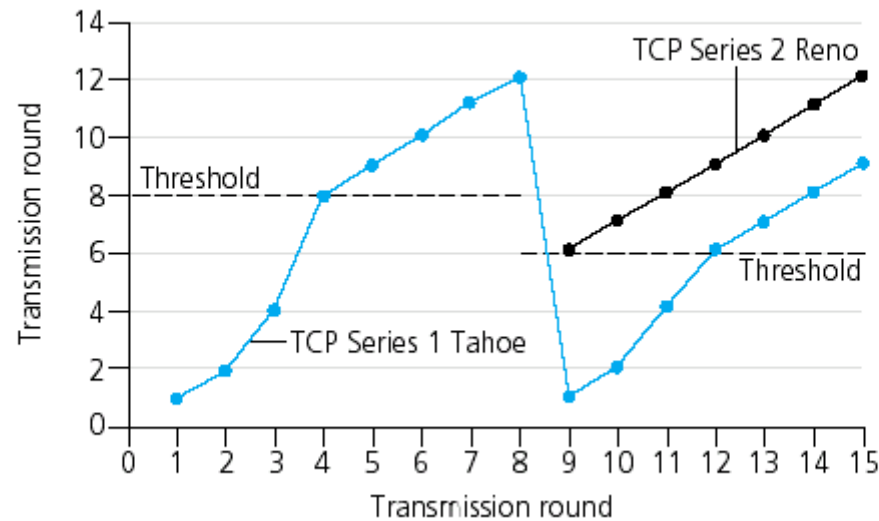❑ timeout indicates a "more alarming" congestion scenario

# Refinement

- Q: When should the exponential increase switch to linear?

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

## Implementation:

- Variable Threshold

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

# Summary: TCP Congestion Control

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

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

- When a triple duplicate ACK occurs, `Threshold` set to `CongWin/2` and `CongWin` set to `Threshold`.

- When timeout occurs, `Threshold` set to `CongWin/2` and `CongWin` is set to 1 MSS.

# TCP sender congestion control

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

# TCP summary

- Connection-oriented: SYN, SYNACK; FIN
- Retransmit lost packets; in-order data: sequence no., ACK no.
- ACKs: either piggybacked, or no-data pure ACK packets if no data travelling in other direction
- Don't overload receiver: rwin
  - rwin advertised by receiver
- Don't overload network: cwin
  - cwin affected by receiving ACKs
- Sender buffer = min { rwin, cwin }
- Congestion control:
  - Slow start: exponential growth of cwin
  - Congestion avoidance: linear groth of cwin
  - Timeout; duplicate ACK: shrink cwin
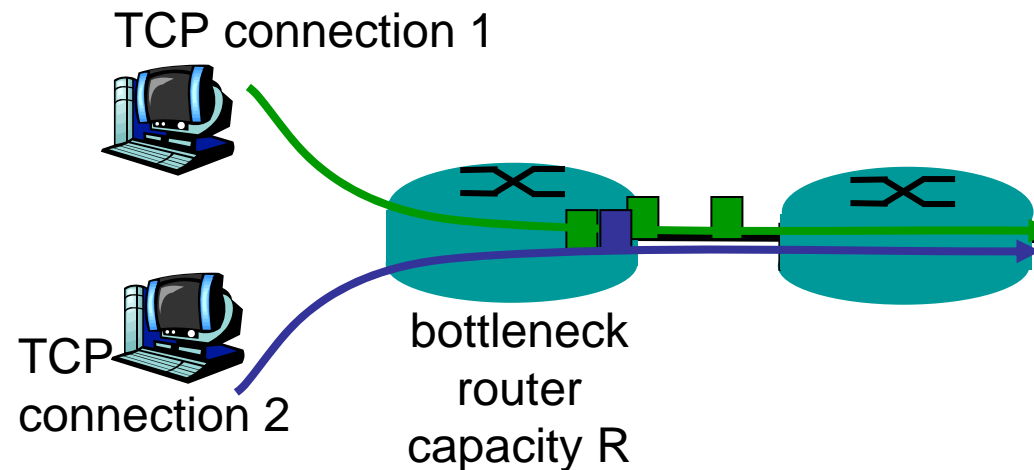- Continuously adjust RTT estimation

# TCP throughput

- ❑ What's the average throughout of TCP as a function of window size and RTT?
  - Ignore slow start
- ❑ Let W be the window size when loss occurs.
- ❑ When window is W, throughput is W/RTT
- ❑ Just after loss, window drops to W/2, throughput to W/2RTT.
- ❑ Average throughout: 0.75 W/RTT

**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
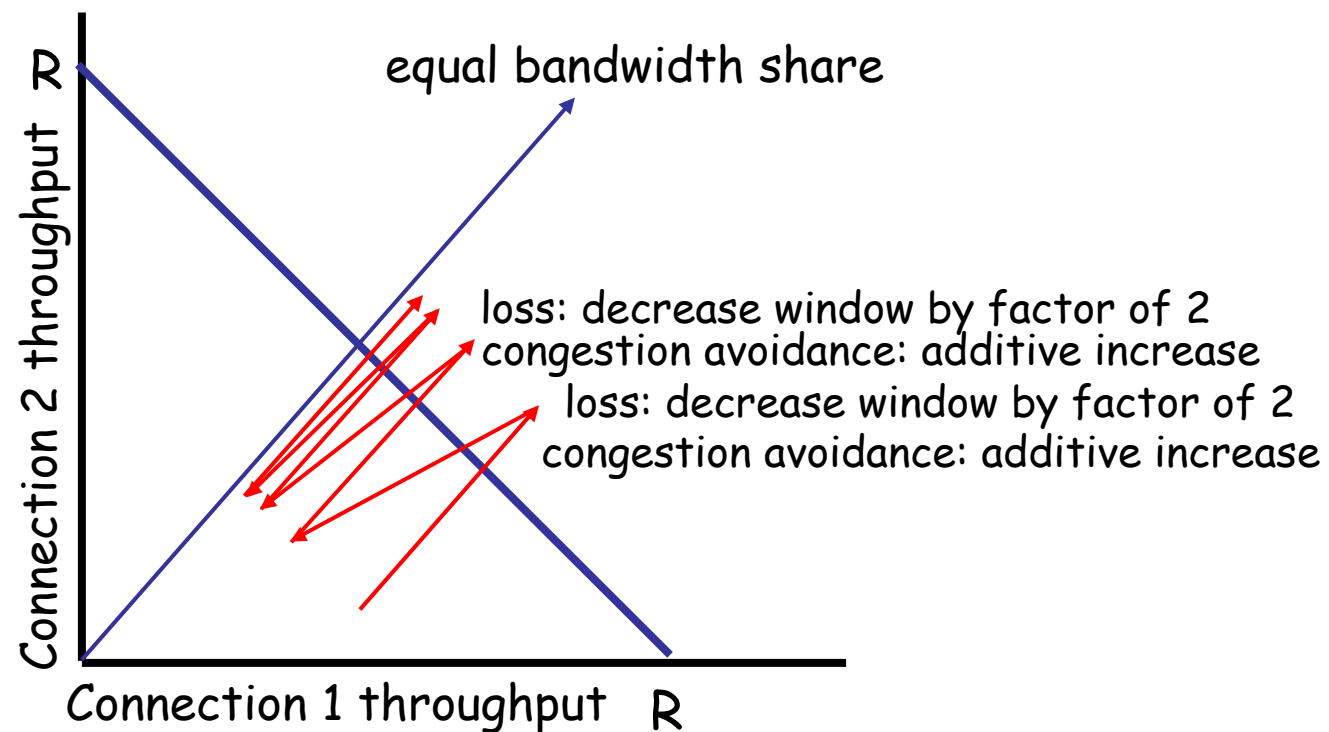
bottleneck router capacity R

## Why is TCP fair?

Two competing sessions:

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

equal bandwidth share

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

Connection 2 throughput

R

Connection 1 throughput R

# Fairness (more)

## Fairness and UDP

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

## Fairness and parallel TCP connections

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

# Chapter 3: Summary

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

<u>Next:</u>

- leaving the network "edge" (application, transport layers)
- into the network "core"