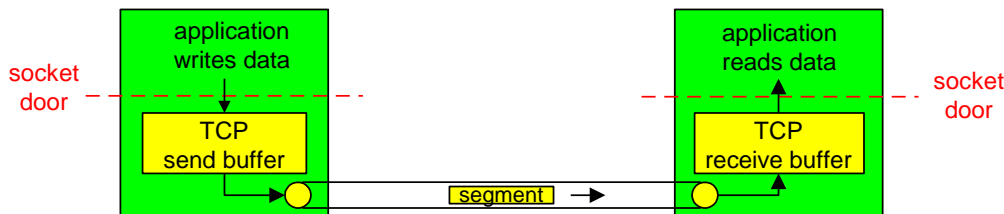# Chapter 3 outline

- 3.1 Transport-layer services

- 3.2 Multiplexing and demultiplexing

- 3.3 Connectionless transport: UDP

- 3.4 Principles of reliable data transfer

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

- 3.6 Principles of congestion control
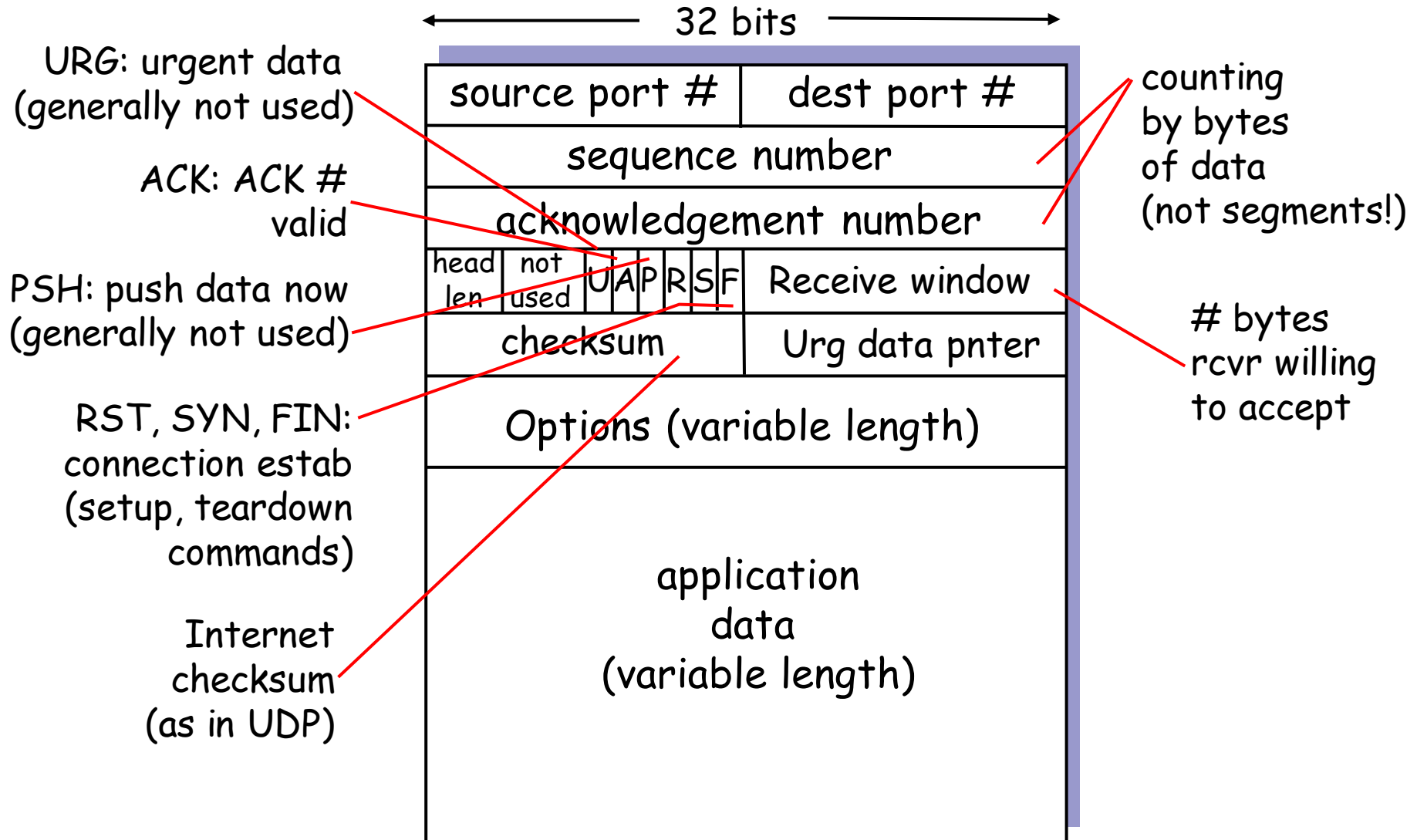
- 3.7 TCP congestion control

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

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

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

socket door

application writes data

TCP send buffer

segment →

application reads data

TCP receive buffer

socket door

# TCP segment structure
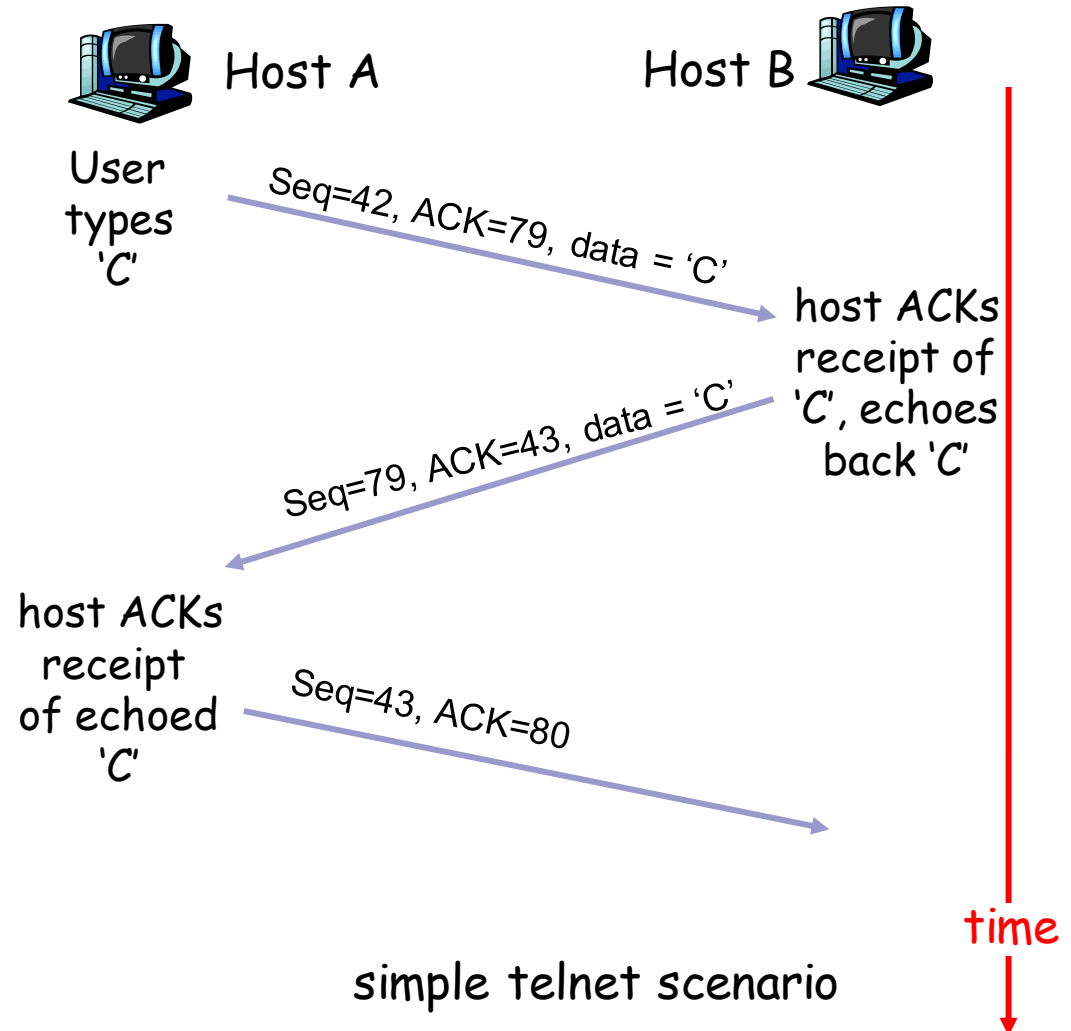
# TCP seq. #'s and ACKs

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

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

Q: how receiver handles out-of-order segments
- □ A: TCP spec doesn't say, - up to implementor

Host A                                    Host B

User types 'C'

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

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

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

host ACKs receipt of echoed 'C'

Seq=43, ACK=80

time

simple telnet scenario

# TCP Round Trip Time and Timeout

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

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

**Q:** how to estimate RTT?

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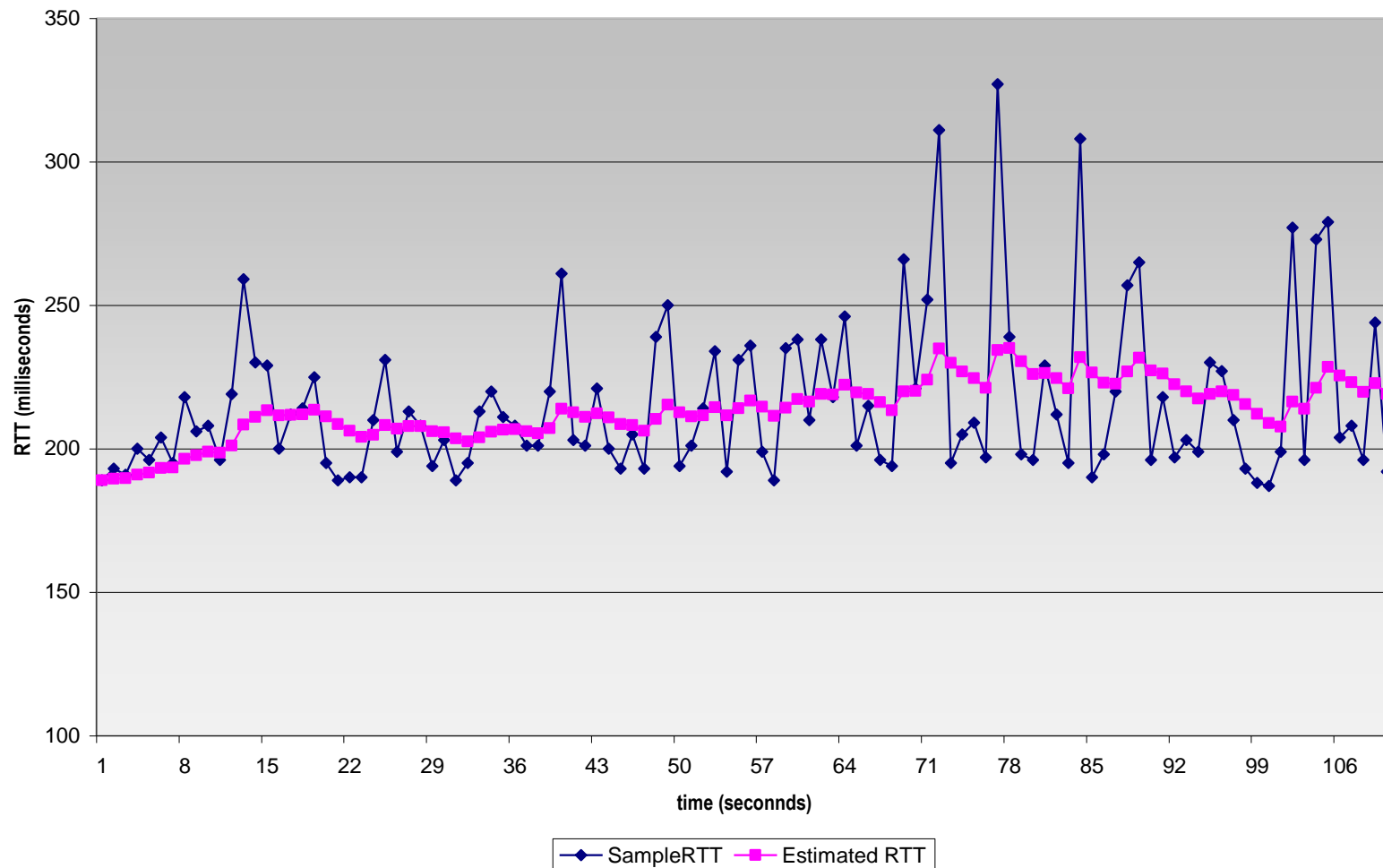  - □ average several recent measurements, not just current **SampleRTT**

# TCP Round Trip Time and Timeout

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

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

# Example RTT estimation:

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

# TCP Round Trip Time and Timeout

## Setting the timeout

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

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


(typically, β = 0.25)
```

Then set timeout interval:

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

# Chapter 3 outline

# 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:

<span style="color:red">data rcvd from app:</span>

- 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`

<span style="color:red">timeout:</span>

- retransmit segment that caused timeout
- restart timer

<span style="color:red">Ack rcvd:</span>

- 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 */
```
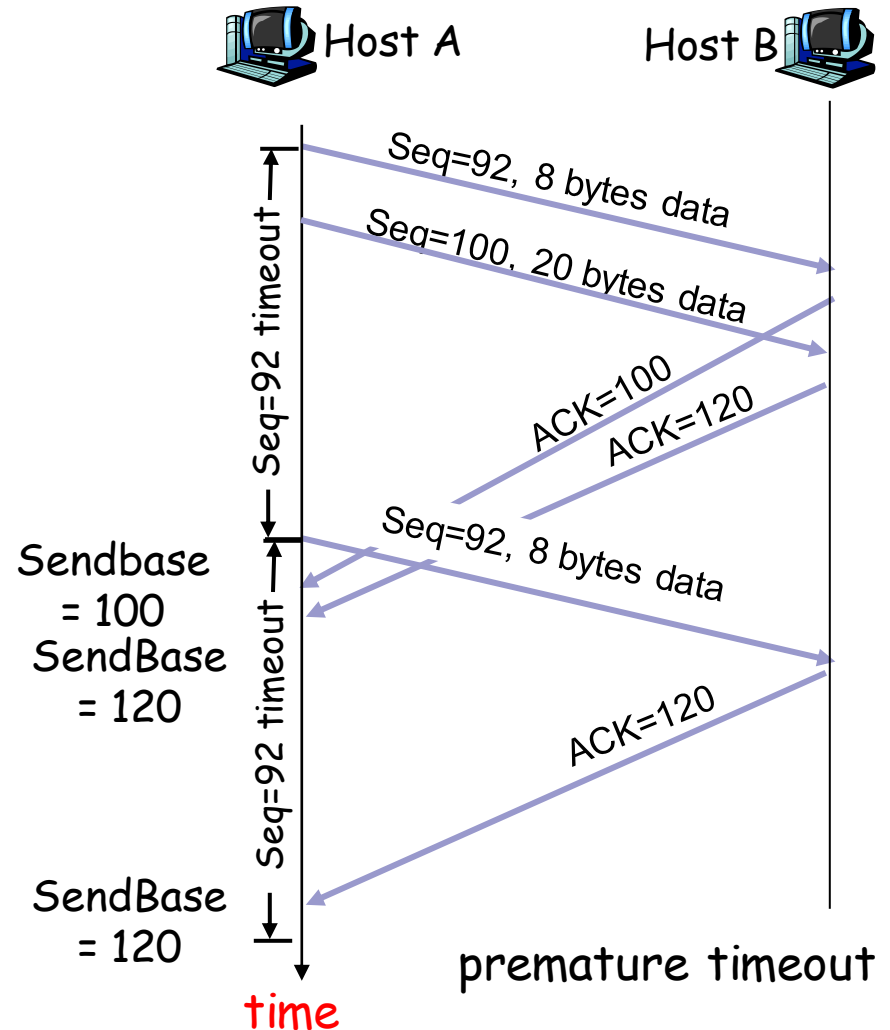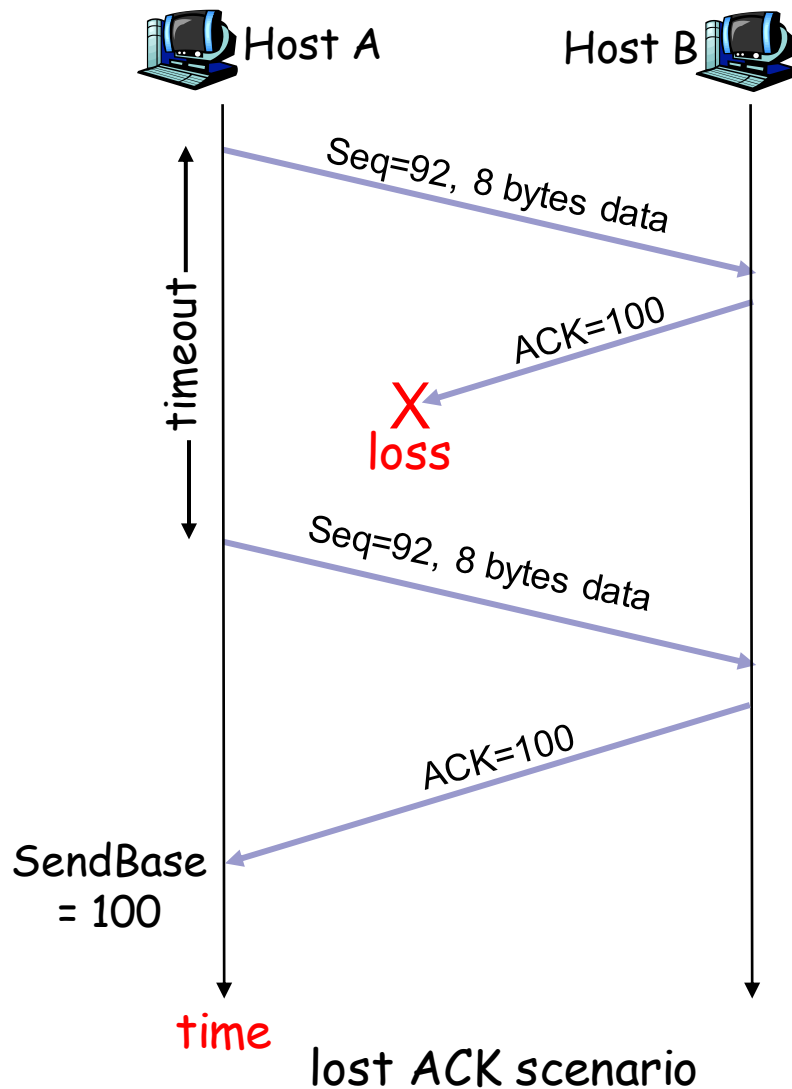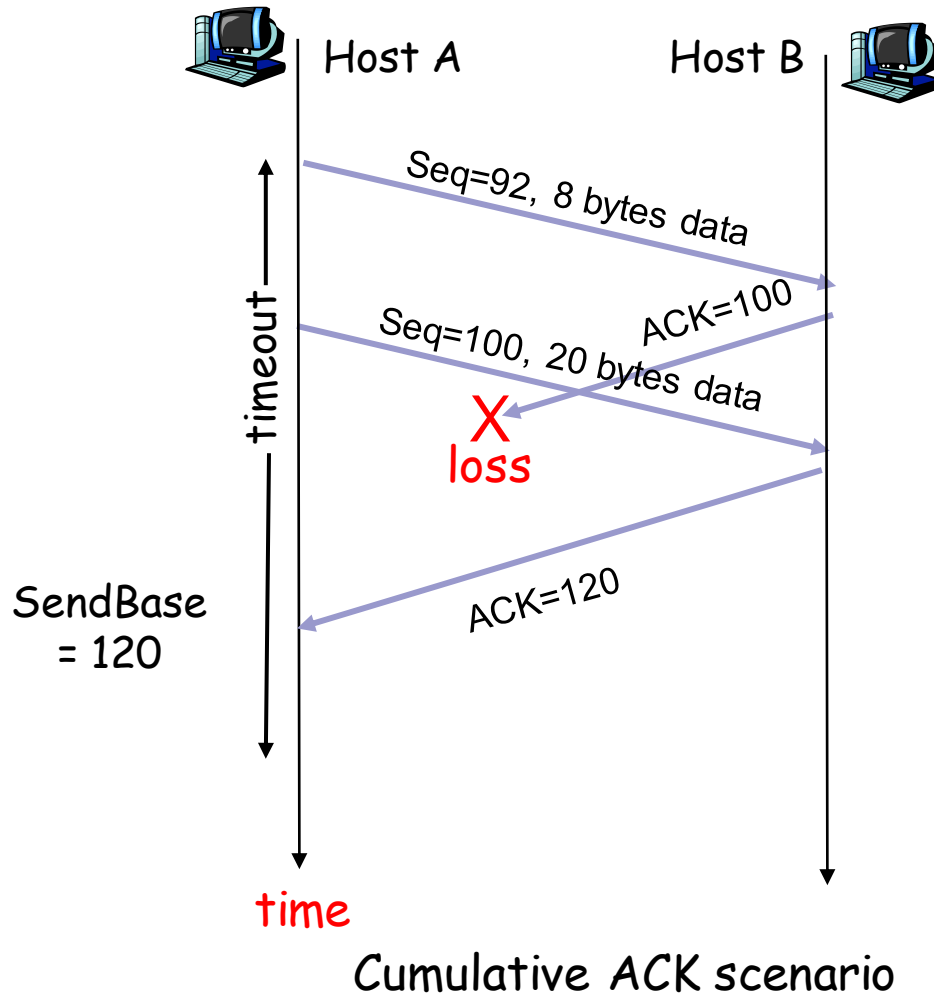
# TCP: retransmission scenarios



lost ACK scenario

premature timeout

# TCP retransmission scenarios (more)



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
            }
```

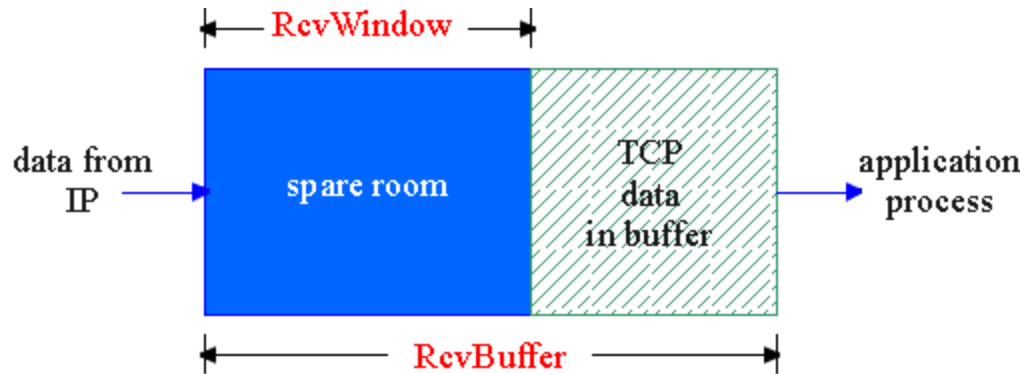a duplicate ACK for
already ACKed segment

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

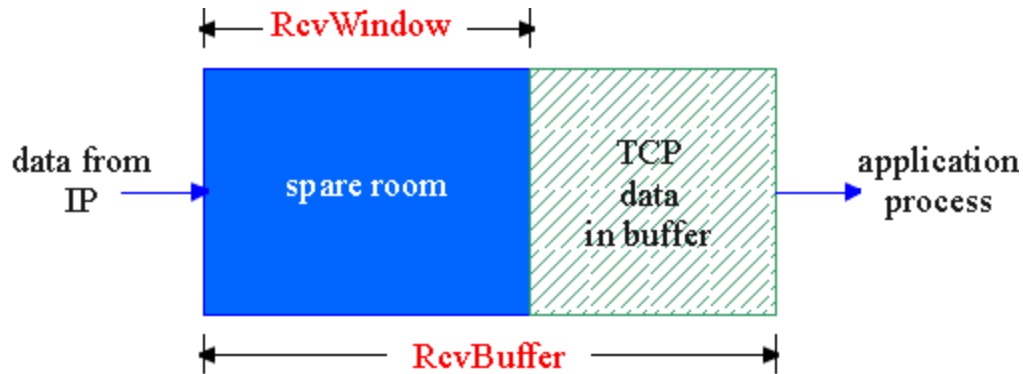- receive side of TCP connection has a receive buffer:

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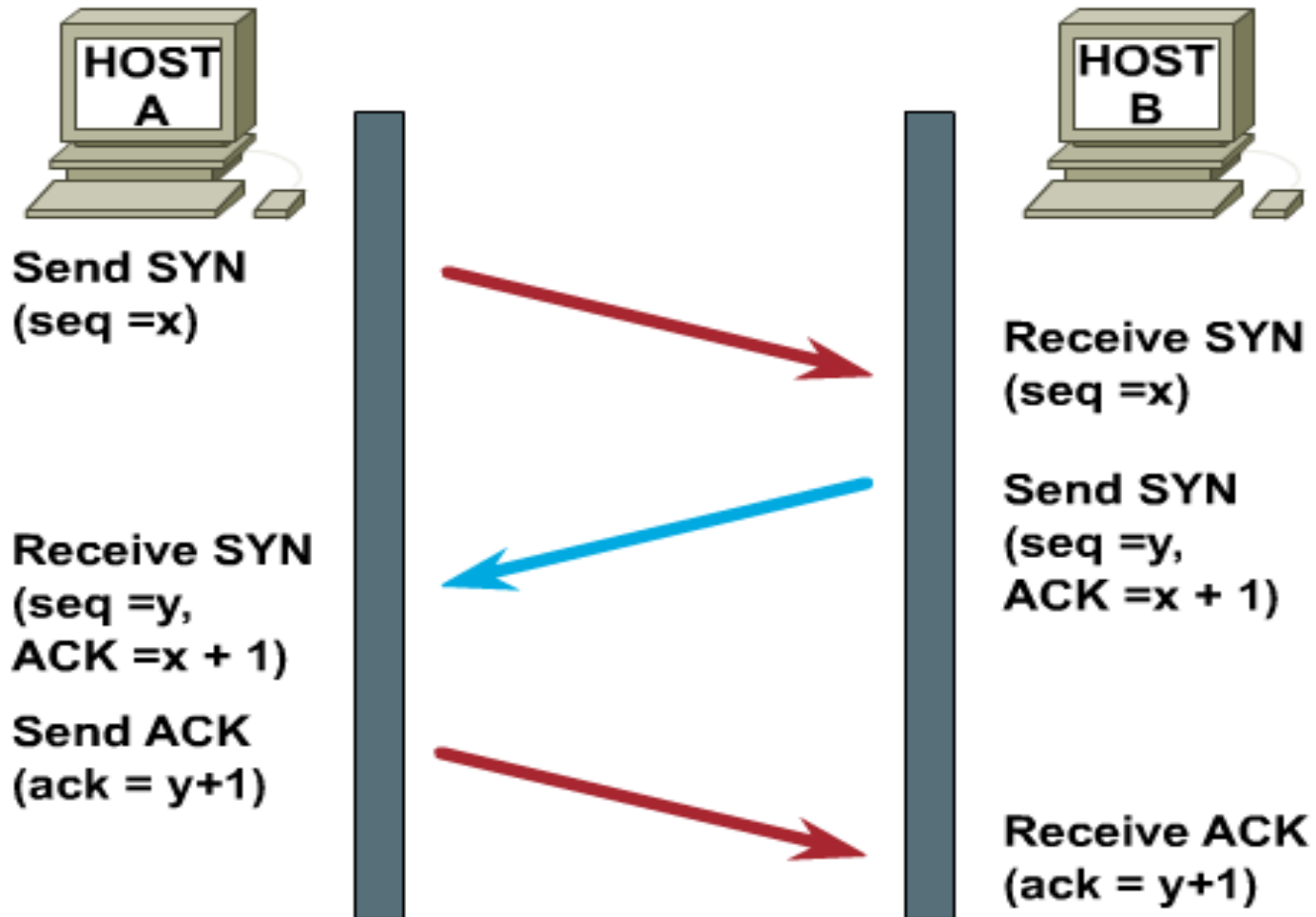
**Step 1:** client host sends TCP SYN segment to server
  - □ specifies initial seq #
  - □ no data

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

  - □ server allocates buffers
  - □ specifies server initial seq. #

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

# Three-Way Handshake



**HOST A**

Send SYN
(seq =x)

Receive SYN
(seq =y,
ACK =x + 1)

Send ACK
(ack = y+1)

**HOST B**

Receive SYN
(seq =x)

Send SYN
(seq =y,
ACK =x + 1)

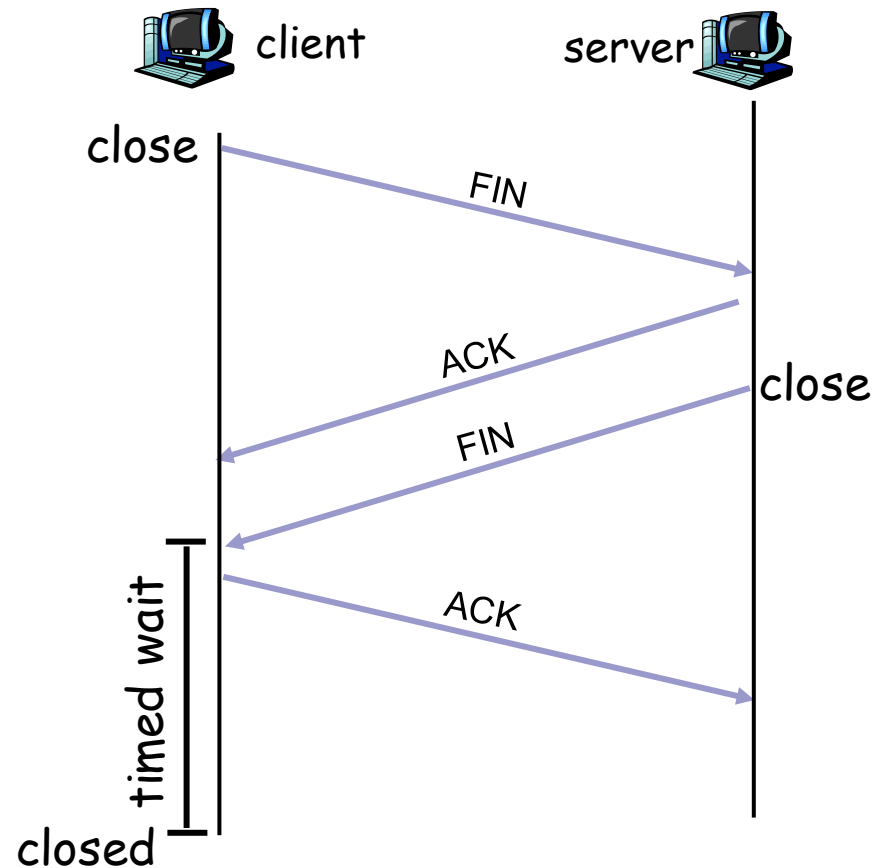Receive ACK
(ack = y+1)

# TCP Connection Management (cont.)

## Closing a connection:

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

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

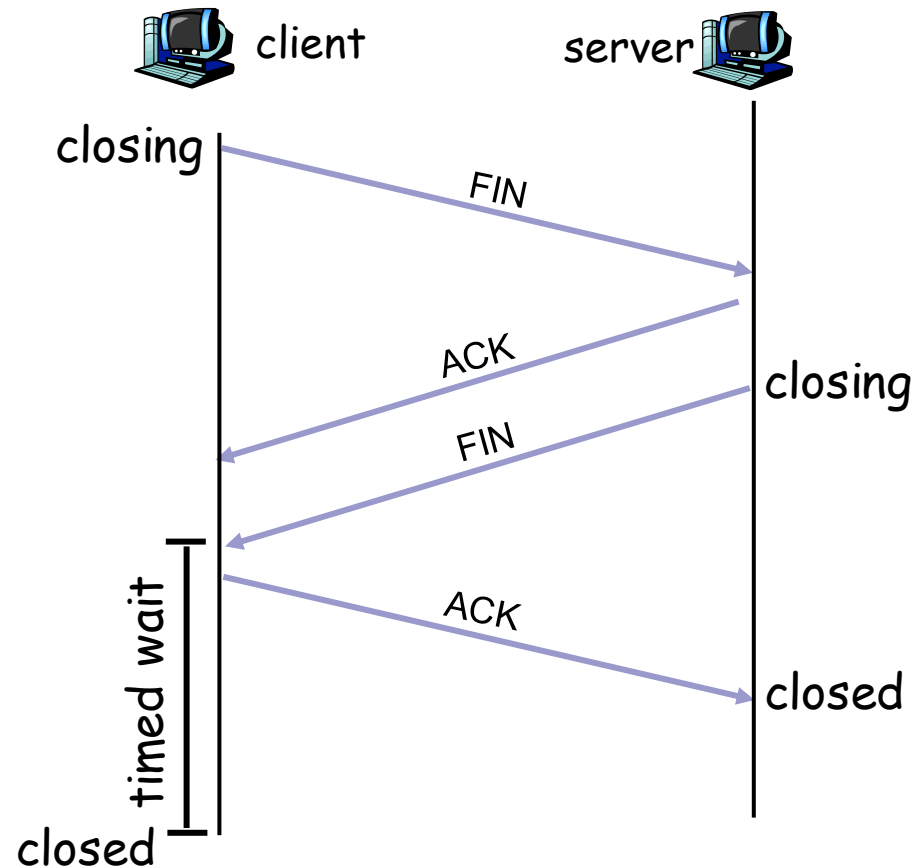Step 2: server receives FIN, replies with ACK. Closes connection, sends FIN.

# TCP Connection Management (cont.)
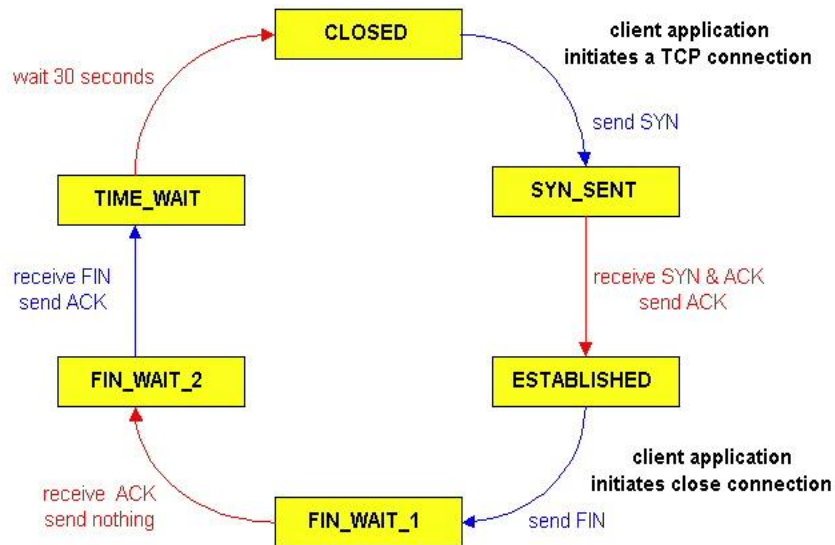
**Step 3:** client receives FIN, replies with ACK.

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

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

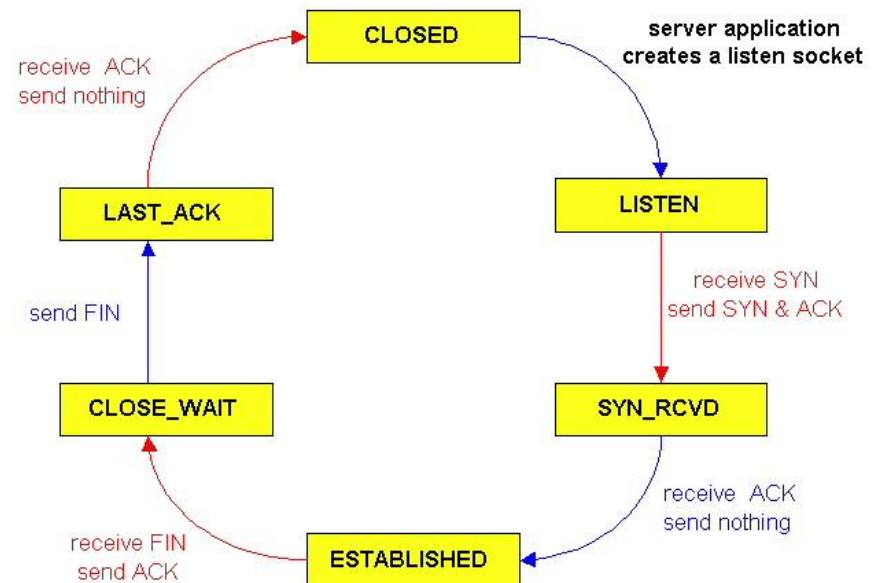**Note:** with small modification, can handle simultaneous FINs.

# TCP Connection Management (cont)



TCP client
lifecycle

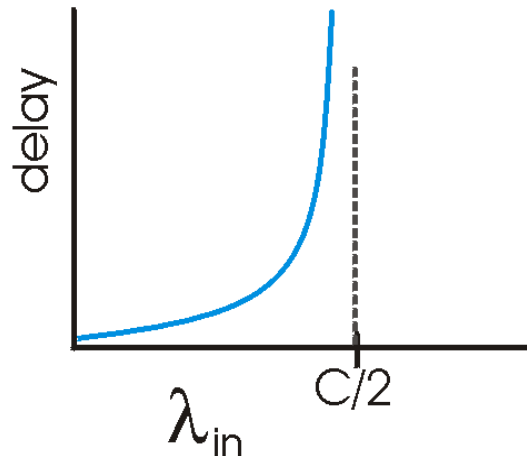TCP server
lifecycle
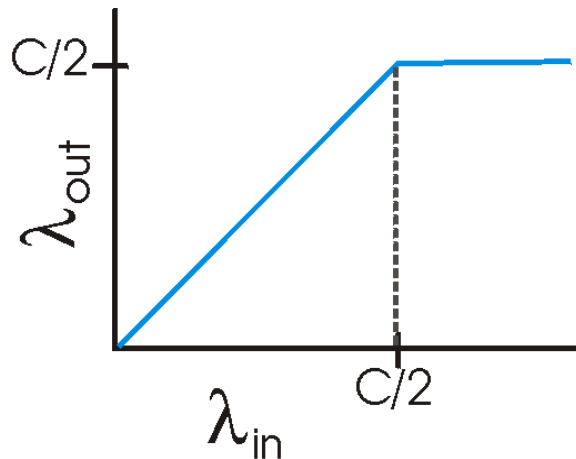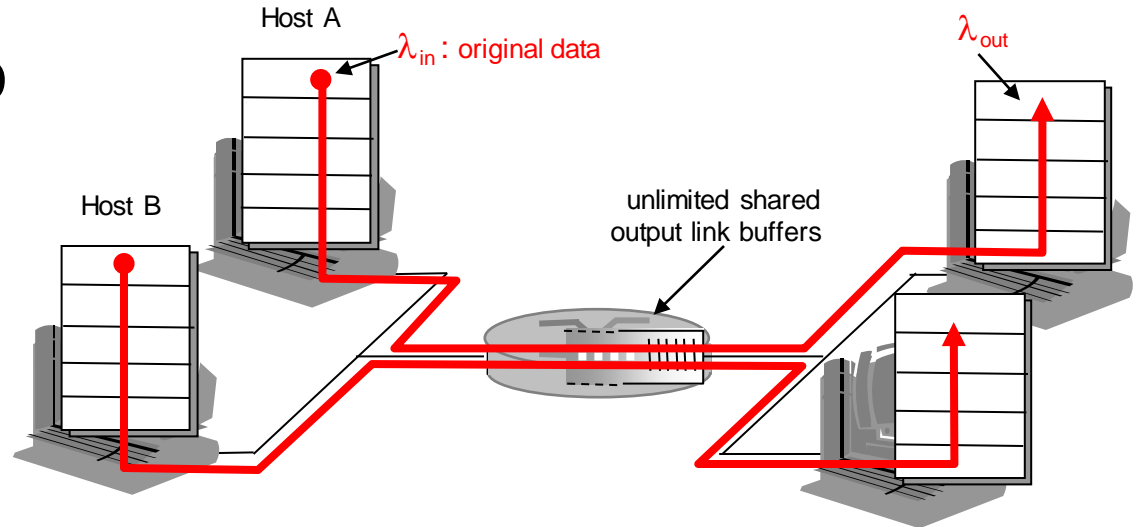
# Chapter 3 outline

# Principles of Congestion Control

Congestion:

- informally: "too many sources sending too much data too fast for *network* to handle"

- different from flow control!

- Manifestations:مظاهر

  - □ lost packets (buffer overflow at routers)

  - □ long delays (queueing in router buffers)

- a top-10 problem!

# Causes/costs of congestion: scenario 1

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



Host A

$\lambda_{in}$ : original data

Host B

unlimited shared output link buffers

$\lambda_{out}$

- large delays when congested
- maximum achievable throughput



$\lambda_{out}$ vs $\lambda_{in}$ with $C/2$



delay vs $\lambda_{in}$ with $C/2$

# Causes/costs of congestion: scenario 2

- one router, *finite* buffers
- sender retransmission of lost packet
- unneeded retransmissions: link carries multiple copies of pkt

Host A

$\lambda_{in}$ : original data

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

$\lambda_{out}$

Host B

finite shared output link buffers

# Causes/costs of congestion



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

# Chapter 3 outline

- 3.1 Transport-layer services
- 3.2 Multiplexing and demultiplexing
- 3.3 Connectionless transport: UDP
- 3.4 Principles of reliable data transfer

- 3.5 Connection-oriented transport: TCP
  - □ segment structure
  - □ reliable data transfer
  - □ flow control
  - □ connection management
- 3.6 Principles of congestion control
- 3.7 TCP congestion control

# TCP Congestion Control

- end-end control (no network assistance)
- sender limits transmission:

  **LastByteSent-LastByteAcked**

  $\leq$ **CongWin**

- Roughly,

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

- **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
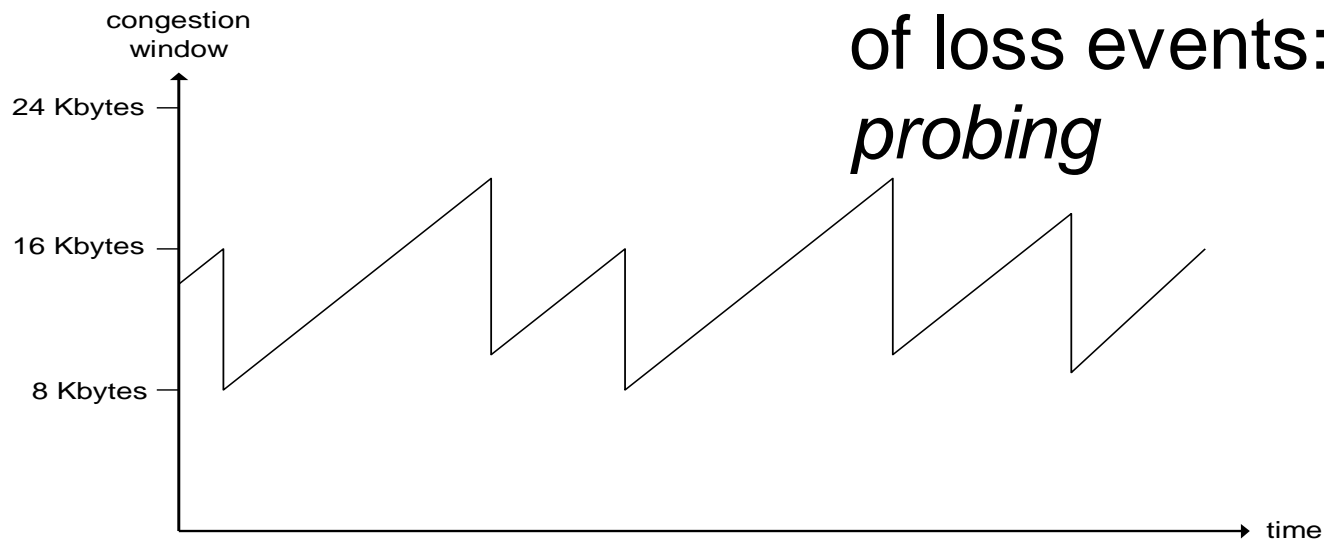- ☐ conservative after timeout events

# TCP AIMD

**multiplicative decrease:** cut `CongWin` in half after loss event

**additive increase:** increase `CongWin` by 1 MSS every RTT in the absence of loss events: *probing*

congestion window

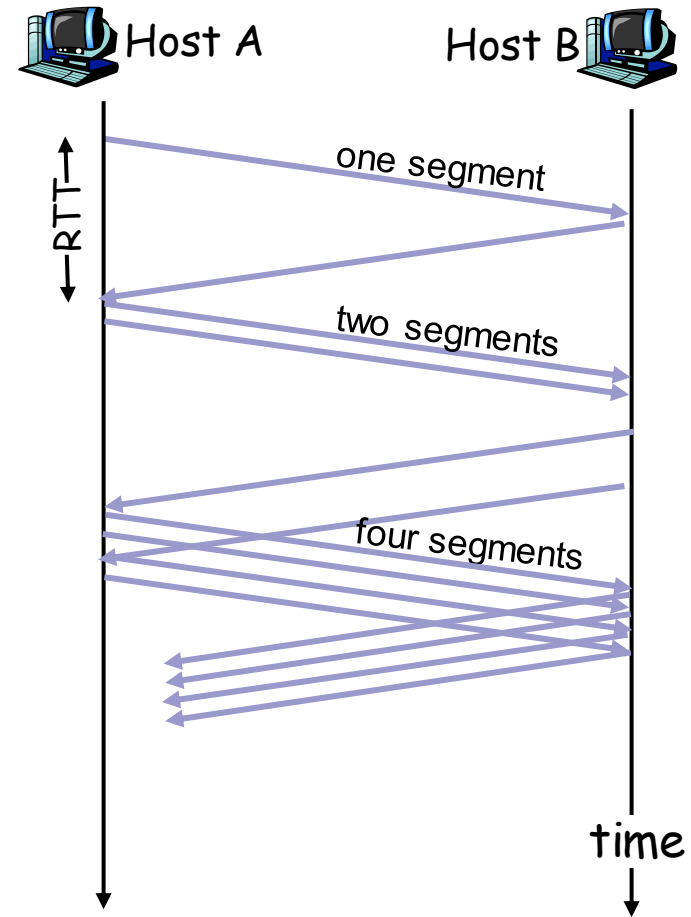24 Kbytes —

16 Kbytes —

8 Kbytes —

time

Long-lived TCP connection

# TCP Slow Start

- 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
- <u>Summary:</u> initial rate is slow but ramps up exponentially fast

# Refinement

- 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
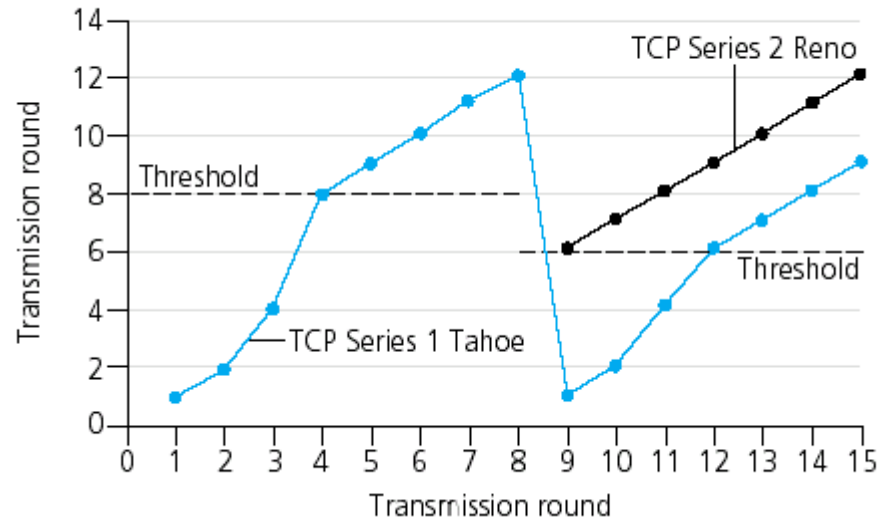• timeout before 3 dup ACKs is "more alarming"

# Refinement (more)

Q: When should the exponential increase switch to linear?

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



## Implementation:

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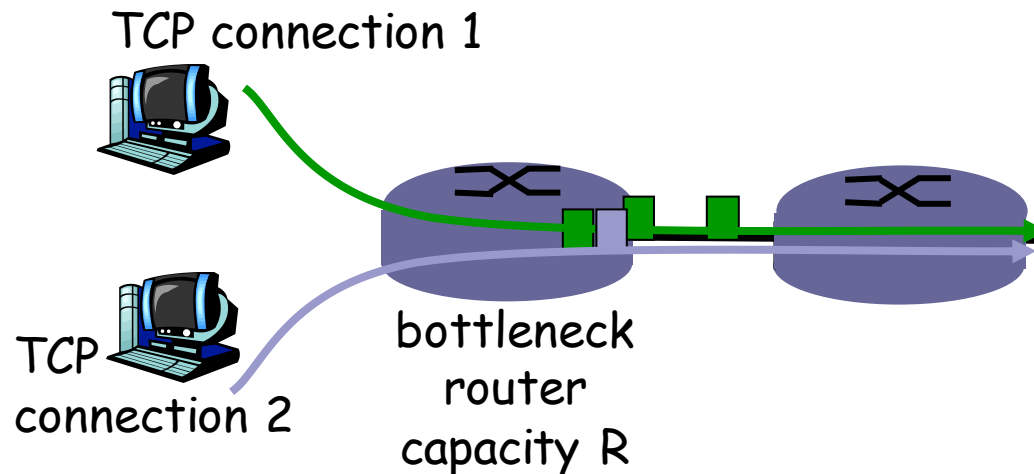- At loss event, Threshold is set to 1/2 of CongWin just before loss event

# 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 Fairness

- Fairness goal: if K TCP sessions share same bottleneck link of bandwidth R, each should have average rate of R/K
- Practically this does not happen in TCP as connections with lower RTT are able to grab the available link bandwidth more quickly.

TCP connection 1

TCP connection 2

bottleneck
router
capacity 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 cnctions between 2 hosts.
- Web browsers do this
- Example: link of rate R supporting 9 cnctions;
  - new app asks for 1 TCP, gets rate R/10
  - new app asks for 11 TCPs, gets R/2 !

# TCP Options: Protection Against Wrap Around Sequence

- 32-bit `SequenceNum`

| Bandwidth | Time Until Wrap Around |
|---|---|
| T1 (1.5 Mbps) | 6.4 hours |
| Ethernet (10 Mbps) | 57 minutes |
| T3 (45 Mbps) | 13 minutes |
| FDDI (100 Mbps) | 6 minutes |
| STS-3 (155 Mbps) | 4 minutes |
| STS-12 (622 Mbps) | 55 seconds |
| STS-24 (1.2 Gbps) | 28 seconds |

# TCP Options: Keeping the Pipe Full

- 16-bit **AdvertisedWindow**

| Bandwidth | Delay x Bandwidth Product |
|---|---|
| T1 (1.5 Mbps) | 18KB |
| Ethernet (10 Mbps) | 122KB |
| T3 (45 Mbps) | 549KB |
| FDDI (100 Mbps) | 1.2MB |
| STS-3 (155 Mbps) | 1.8MB |
| STS-12 (622 Mbps) | 7.4MB |
| STS-24 (1.2 Gbps) | 14.8MB |

assuming 100ms RTT

# TCP Extensions

- Implemented as header options
- Store timestamp in outgoing segments
- Extend sequence space with 32-bit timestamp (PAWS)
- Shift (scale) advertised window