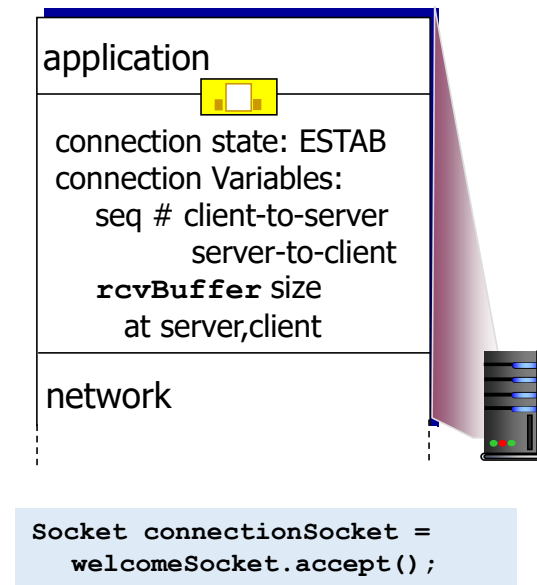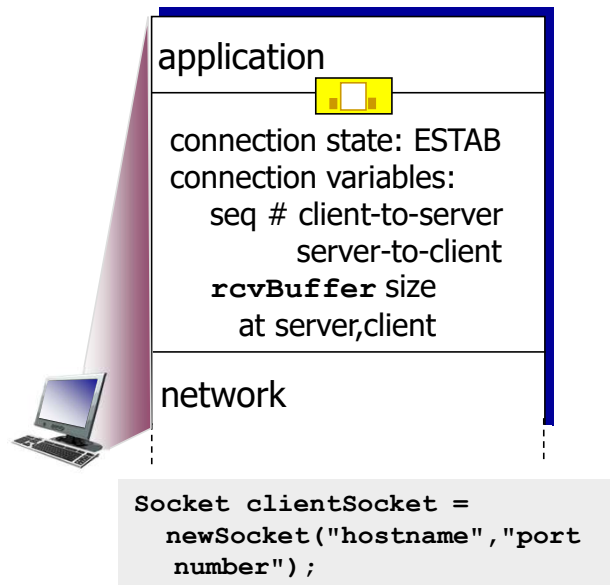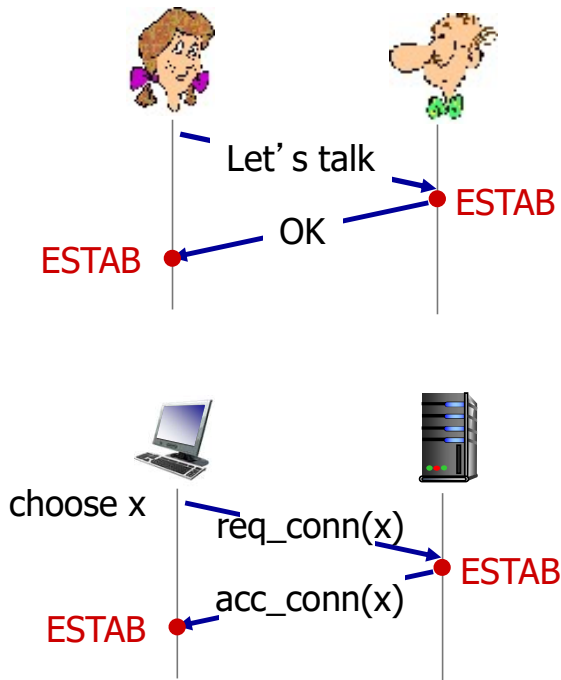# CONNECTION MANAGEMENT

before exchanging data, sender/receiver "handshake":

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



application

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

network

```
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();
```

## 2-way handshake:

choose x

**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

Problems with 2-way handshake:

i) Early/Premature connection establishment -- false connections.

ii) Half-open connections – Duplicate packets from previous connection.

iii) Spoofing Problem – Bombard server with loads of unrelated data

*client state*

*server state*

```
cSocket = socket(AF_INET,SOCK_STREAM)
```

```
sSocket = socket(AF_INET,SOCK_STREAM)
sSocket.bind(('',sPort))
sSocket.listen(1)
connectionSocket, Caddr = sSocket.accept()
```

LISTEN

```
cSocket.connect((sName,sPort))
```

choose init
seq num, x
send TCP SYN
msg

SYNbit=1, Seq=x

LISTEN

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

SYNSENT

SYNbit=1, Seq=y
ACKbit=1; ACKnum=x+1

SYN RCVD

received SYNACK(x)
indicates server is live;
send ACK for SYNACK;
**this segment may
contain
client-to-server
data**

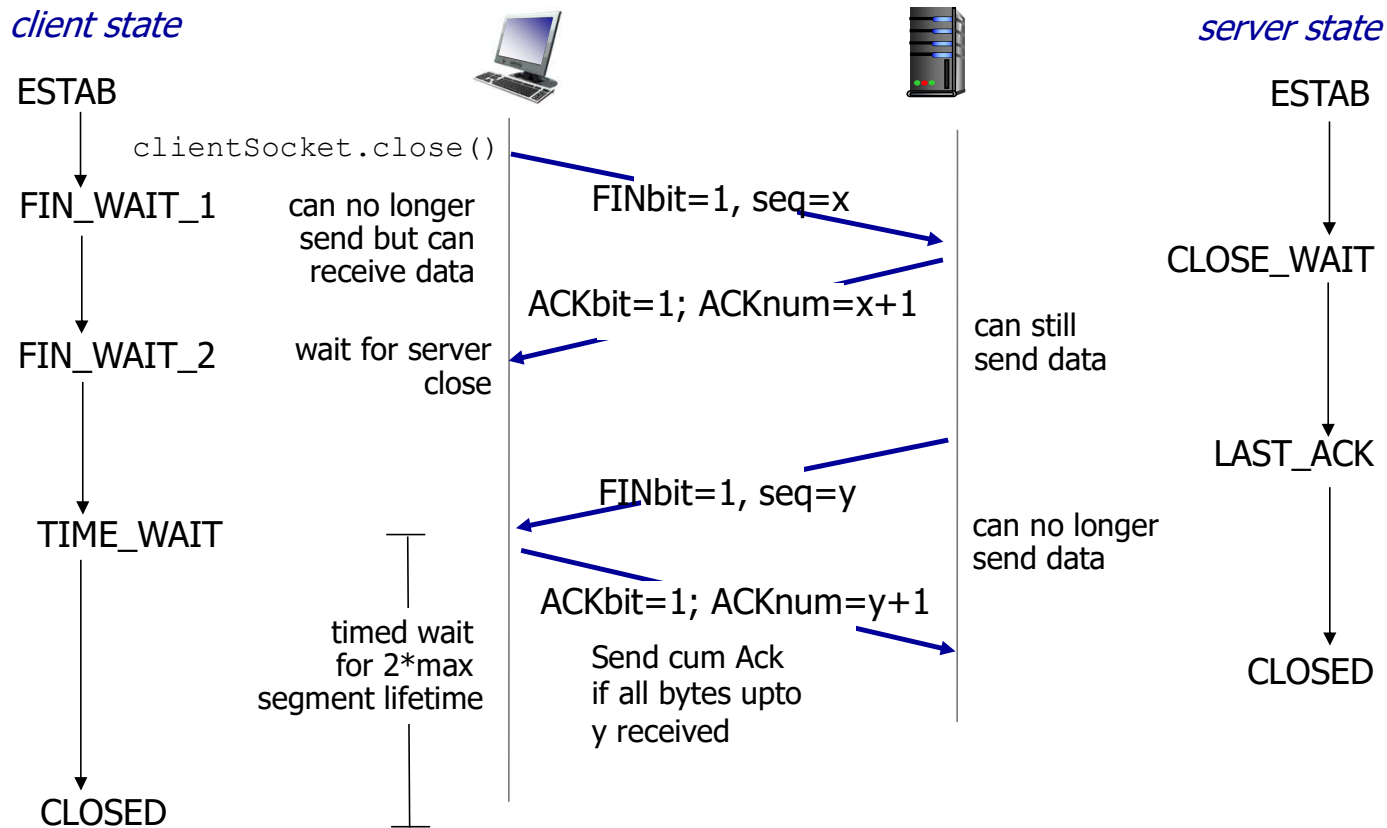ACKbit=1, ACKnum=y+1

ESTAB

received ACK(y)
indicates client is live

ESTAB

RFC 793: The principle reason for the three-way handshake to reduce the possibility of false connections i.e. to prevent old duplicate connection initiations from causing confusion.

+

Bi-Directional ISN Synchronization

client state

server state

ESTAB

ESTAB

`clientSocket.close()`

FIN_WAIT_1

can no longer
send but can
receive data

FINbit=1, seq=x

CLOSE_WAIT

ACKbit=1; ACKnum=x+1

can still
send data

FIN_WAIT_2

wait for server
close

FINbit=1, seq=y

LAST_ACK

TIME_WAIT

can no longer
send data

ACKbit=1; ACKnum=y+1

timed wait
for 2*max
segment lifetime

Send cum Ack
if all bytes upto
y received

CLOSED

CLOSED

- Important to pick ISN appropriately; Choosing a fixed ISN is undesirable.

- Pick random value (intended to avoid others predicting ISN for a new connection)

**sequence numbers:**
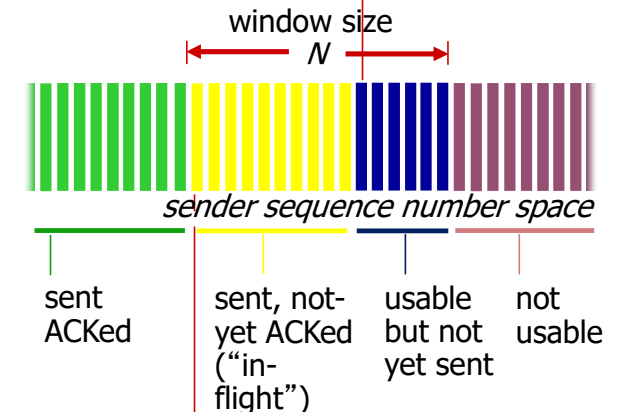- byte stream "number" of the first byte in segment's data

**acknowledgement number:**
- seq # of **next** byte expected from other side
- Support cumulative ACK

**Q:** how should the receiver handle out-of-order segments?
- A: TCP spec doesn't say, - it is up to the implementor..

outgoing segment from sender

| source port # | dest port # |
|---|---|
| sequence number | |
| acknowledgement number | |
| | rwnd |
| checksum | urg pointer |

window size
N

*sender sequence number space*

| sent ACKed | sent, not-yet ACKed ("in-flight") | usable but not yet sent | not usable |

incoming segment to sender

| source port # | dest port # |
|---|---|
| sequence number | |
| acknowledgement number | |
| A | rwnd |
| checksum | urg pointer |

## Q: How do you detect loss of packets ?
## *And*
## *How do you recover from packet loss?*

Packet loss scenario

lost ACK scenario

cumulative packets with lost ACK

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

❖longer than RTT - but RTT varies over time.

❖*too short:* premature timeout, unnecessary retransmissions.

❖*too long:* slow reaction to segment loss.

**Q:** how to estimate RTT?

- **SampleRTT(R)**: measured time from segment transmission until ACK receipt.
  - SampleRTT will vary, we want estimated RTT to be "smoother" using EWMA ($\alpha = 0.125$)
  - average several *recent* measurements, not just use current **SampleRTT**
  - ignore retransmissions

**RTO**: Retransmission Timeout – duration to wait before retransmission in the absence of any ack for data.

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



*SampleRTT = $RTT_i$; measured for every sample (#Seq <-> ACK pair)*

**Initialization**: Until SampleRTT is measured set *RTO =1s.*

When the first RTT measurement $R_0$ is made, the host MUST set

$SRTT \leftarrow R_0$

$RTTVAR \leftarrow R_0/2$

$RTO \leftarrow SRTT + MAX (G, K*RTTVAR)$, where $K = 4$ and $G$ is the clock granularity.

- For subsequent samples of RTT ($R_i$) where i>=1:

- First estimate SampleRTT *deviation (RTTVAR)* from current SampleRTT and past SmoothedRTT;

- Then estimate SmoothedRTT (SRTT)

$$RTTVAR_{(i)} = (1-\beta)*RTTVAR_{(i-1)} + \beta*|R_{(i)}-SRTT_{(i-1)}|$$
$$SRTT_{(i)} = (1-\alpha)*SRTT_{(i-1)} + \alpha*R_{(i)}$$

(typically, $\beta = 0.25$ and $\alpha = 0.125$)

- **timeout interval: SmoothedRTT** plus "safety margin"

  - large variation in **SmoothedRTT ->** larger safety margin

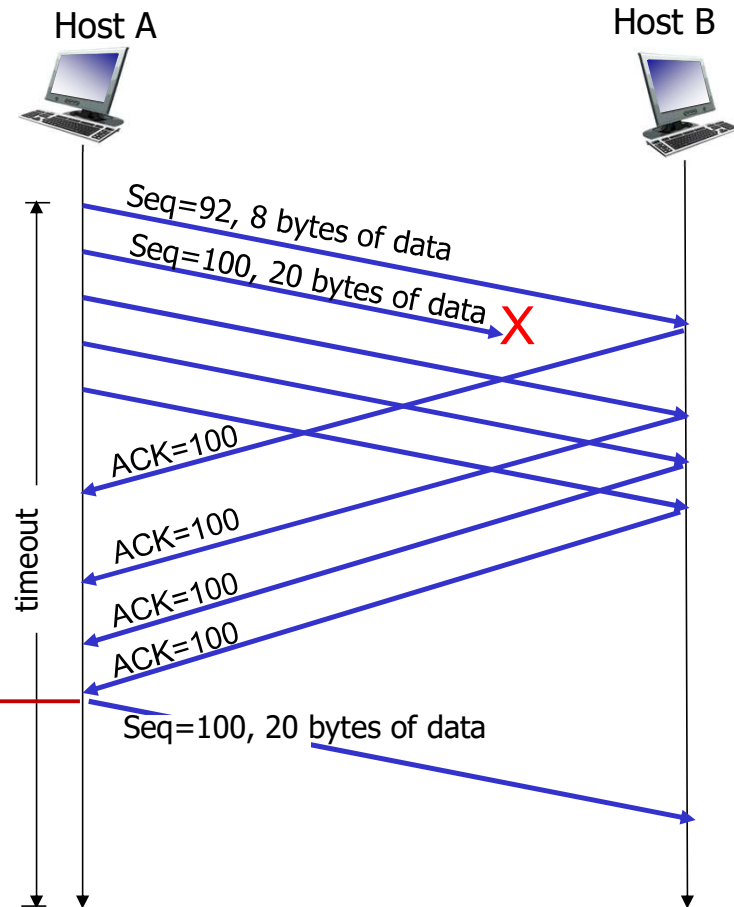**TimeoutInterval (RTO) = *SmoothedRTT* + 4*DevRTT**

estimated RTT          "safety margin"

## TCP fast retransmit

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

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

💡 Receipt of three duplicate ACKs indicates 3 segments received after a missing segment – lost segment is likely. So retransmit!
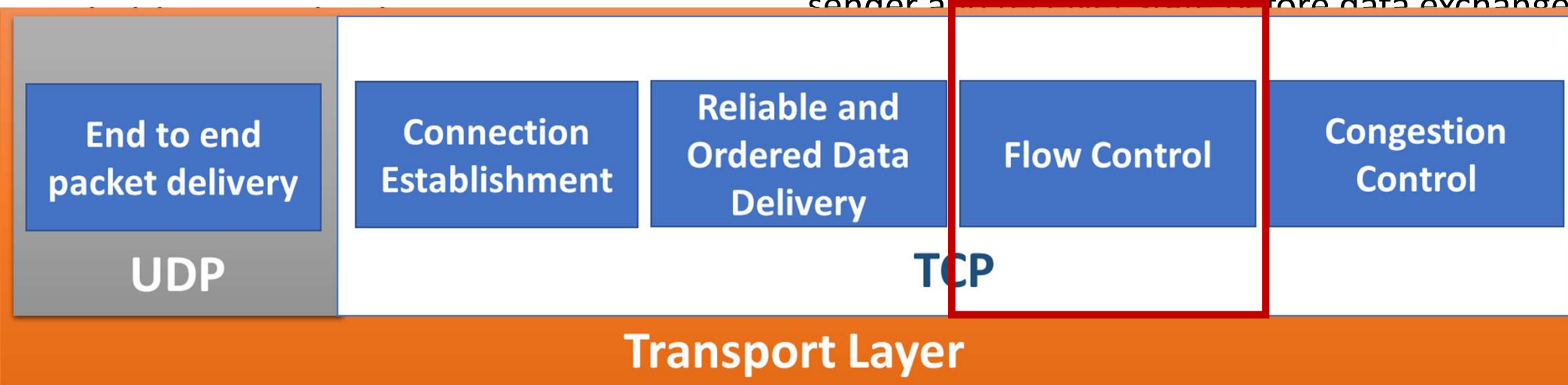
Host A

Host B

Seq=92, 8 bytes of data

Seq=100, 20 bytes of data ✗

ACK=100

ACK=100

ACK=100

ACK=100

Seq=100, 20 bytes of data

timeout

- **point-to-point:**
  - one sender, one receiver

- **connection-oriented:**
  - handshaking (exchange of control msgs) inits
    sender and receiver state before data exchange

| End to end packet delivery | Connection Establishment | Reliable and Ordered Data Delivery | Flow Control | Congestion Control |
|---|---|---|---|---|

**UDP**                                           **TCP**
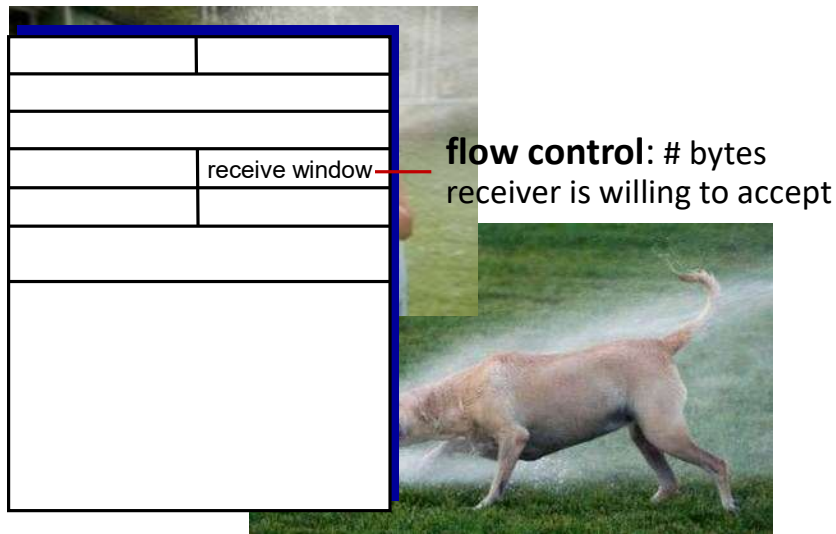
**Transport Layer**

- TCP congestion and flow control set
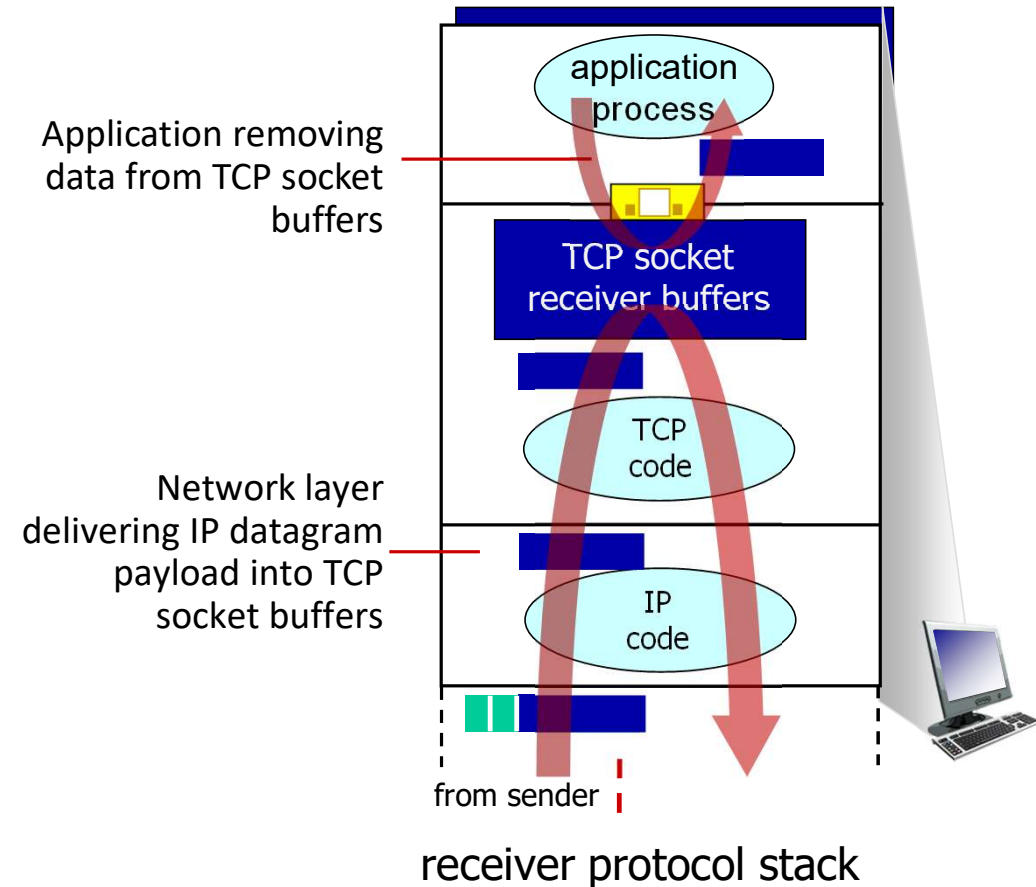  window size

to the network capacity.

Today's Focus:
Flow Control
Congestion Control

# TCP FLOW CONTROL

*Q:* What happens if network layer delivers data faster than application layer removes data from socket buffers?

**flow control**: # bytes receiver is willing to accept

receive window

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

Application removing data from TCP socket buffers

Network layer delivering IP datagram payload into TCP socket buffers

application process

TCP socket receiver buffers

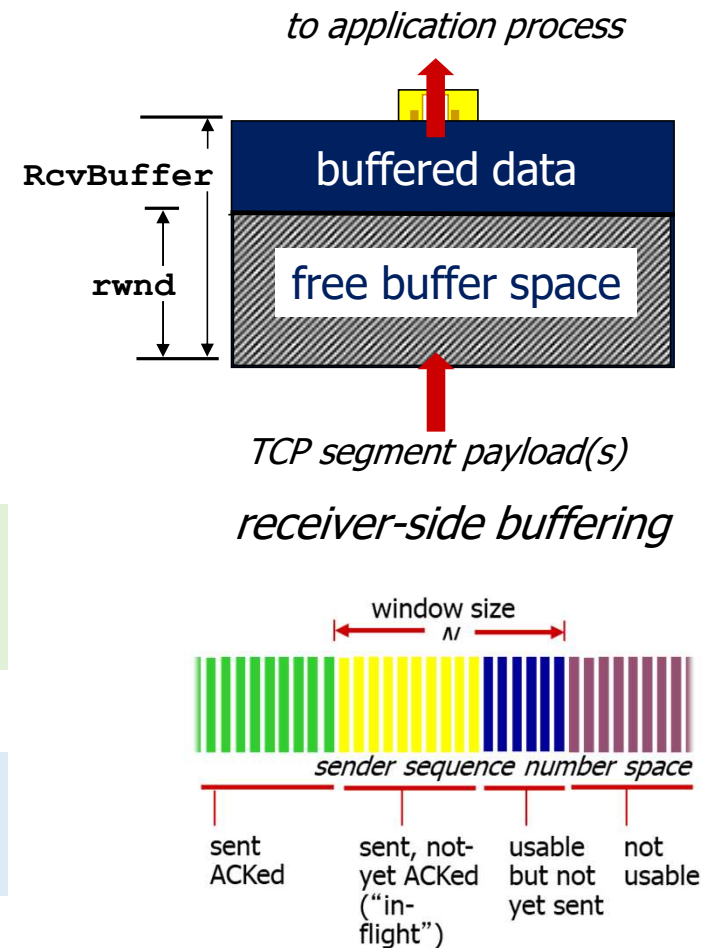TCP code

IP code

from sender

receiver protocol stack

- receiver "advertises" free buffer space by including **rwnd** value in TCP header of receiver-to-sender segments
  - **RcvBuffer** size set via socket options (default 64K bytes)
  - many operating systems auto adjust **RcvBuffer**

- sender limits amount of unacked ("in-flight") data to receiver's **rwnd** value

- guarantees receive buffer will not overflow

Receiver:

rwnd = RcvBuffer – [LastByteRcvd – LastByteRead]

Sender:

LastByteSent – LastByteAcked <= rwnd

*to application process*

*RcvBuffer*

buffered data

*rwnd*

free buffer space

*TCP segment payload(s)*

*receiver-side buffering*

window size
*N*

*sender sequence number space*

sent
ACKed

sent, not-
yet ACKed
("in-
flight")

usable
but not
yet sent

not
usable

- Advancing a full window
  - When the receive window fills up, how to get started again?
    - Sender sends periodic probe to check availability while recv. window=0

- 'Silly Window Syndrome': Very small sized window ➔ small data transfer.
  - receiver requests sender to reduce the amount of data sent at a time.
  - Ack. with-holding at the receiver (delayed acks) helps, but not complete solution.

- 'Tinygram Syndrome': small data at the sender ➔ small data transfer.
  - Application at the sender has very small amount of data to be sent at a time.
  - Coalesce multiple chunks of small data into large segment.
  - Nagle's algorithm: RFC-896
    - If send.window < 1 MSS, delay sending the packet and aggregate/batch more data
    - Can be over-ridden by TCP-NODELAY option.

- **point-to-point:**
  - one sender, one receiver

- **connection-oriented:**
  - handshaking (exchange of control msgs) inits sender and receiver state before data exchange



- TCP congestion and flow control set window size

to the network capacity.

Today's Focus:
Congestion Control

*congestion*:

- informally:

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

  **Deadlock and/or a Livelock**



- manifestations:
  - lost packets (buffer overflow at routers)
  - long delays (queueing in router buffers)
- a top-10 networking problem!

Congestion – we are familiar with…

*congestion*: *<Overload on Network>*

- application-layer input: $\lambda_{in}$ = application-layer output $\lambda_{out}$
- Sender retransmits lost, timed-out packet(s).
- transport-layer input includes *retransmissions* : $\lambda'_{in} \geq \lambda_{in}$

$\lambda_{in}$ : original data

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

$\lambda_{out}$

Host A

SVC-1

Host B

finite link buffers/ shared output

SVC-2

**"costs" of congestion:**

- Large delays (more queue ➔ more delay)
- more work (lost packets ➔ retransmissions)
- lower "goodput" or achieved receiver throughput.
- duplicates: link carries multiple copies of a packet.

**Q1:** What causes congestion?

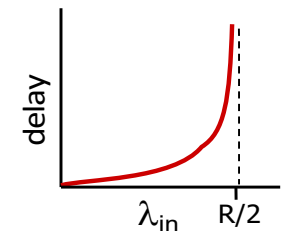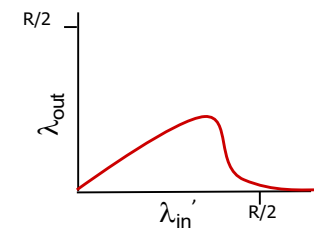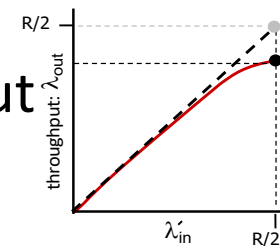throughput: $\lambda_{out}$ — R
$\lambda_{in}$ — R

delay
$\lambda_{in}$ — R

**Q2:** What happens as more hosts join?

throughput: $\lambda_{out}$ — R/2
$\lambda_{in}$ — R/2

delay
$\lambda_{in}$ — R/2
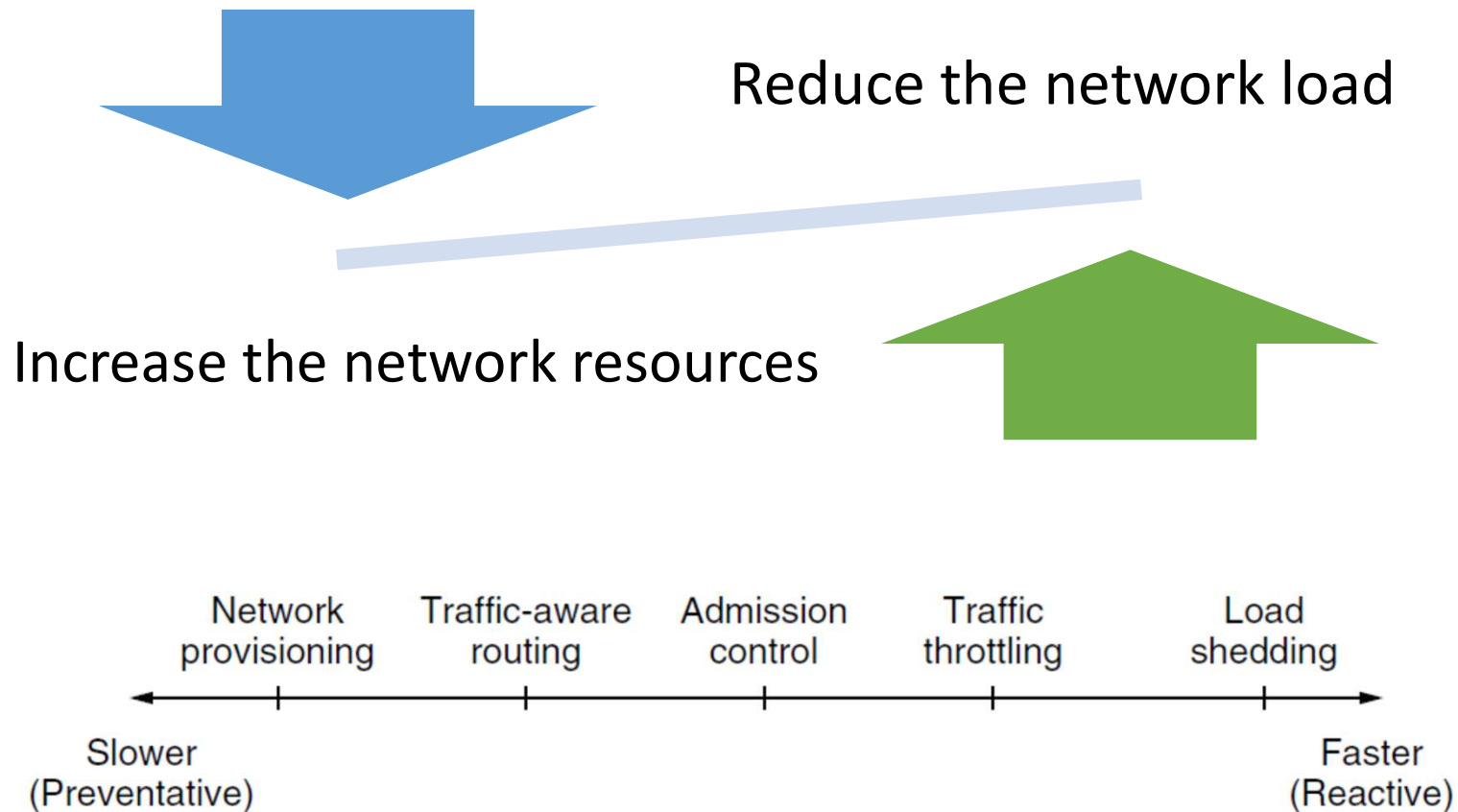
**Q3:** What is/are the cost(s) of congestion?

- throughput can never exceed capacity

- delay increases as capacity approached

- loss/retransmission decreases effective throughput

- un-needed duplicates further decreases effective throughput

- upstream transmission capacity / buffering wasted for packets lost downstream

Reduce the network load

Increase the network resources
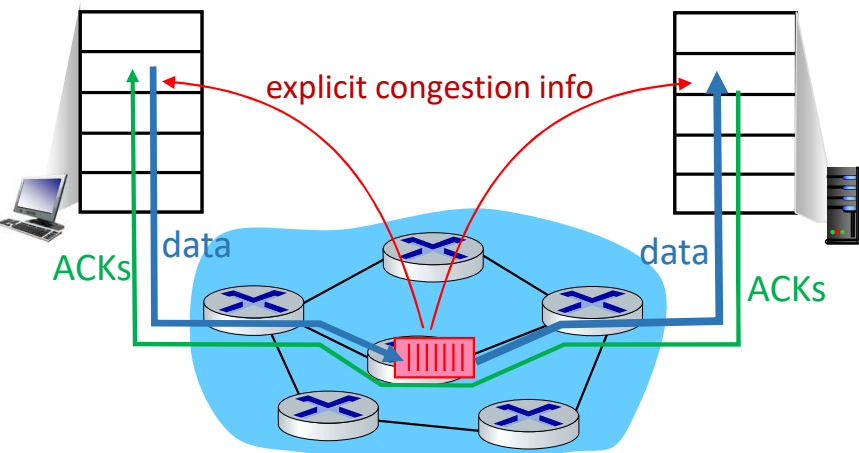
| Network provisioning | Traffic-aware routing | Admission control | Traffic throttling | Load shedding |

Slower (Preventative) ← → Faster (Reactive)

**two broad approaches towards congestion control:**
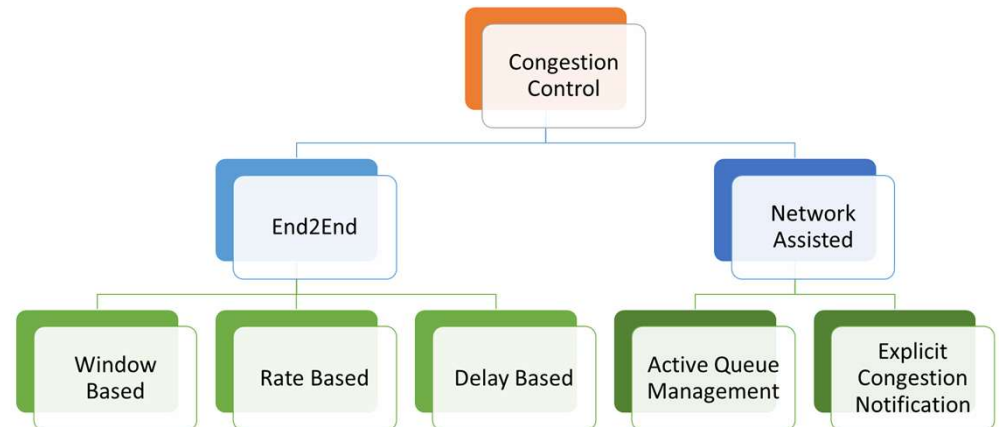
## end-end congestion control:

❖ congestion inferred from end-system observed loss, delay.

❖ no explicit feedback from network
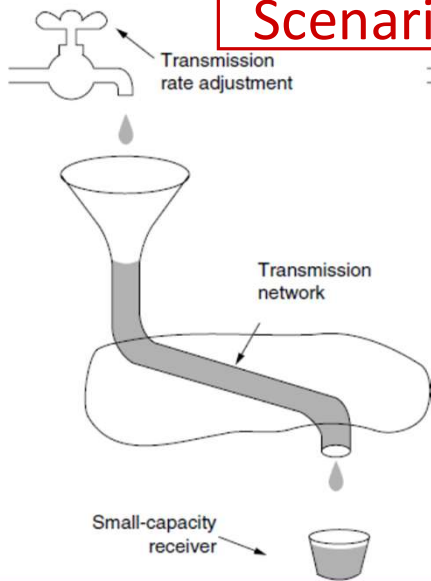
explicit congestion info

ACKs    data    data    ACKs

## network-assisted congestion control:

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

❖ Sender and receiver react to the feedback

Congestion Control

End2End          Network Assisted

Window Based    Rate Based    Delay Based    Active Queue Management    Explicit Congestion Notification
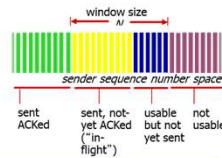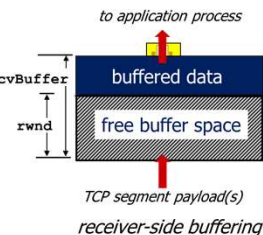
## Scenario 1: A familiar one



**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 (default 64K bytes)
  - many operating systems auto adjust **RcvBuffer**
- sender limits amount of unacked ("in-flight") data to receiver's **rwnd** value
- guarantees receive buffer will not overflow

to application process

buffered data

free buffer space

TCP segment payload(s)

*receiver-side buffering*

Receiver:
rwnd = RcvBuffer – [LastByteRcvd – LastByteRead]

Sender:
LastByteSent – LastByteAcked <= rwnd

*sender sequence number space*

sent ACKed | sent, not-yet ACKed ("in-flight") | usable but not yet sent | not usable

## Scenario 2: Network Congestion

❖ **cwnd** is dynamic, function of perceived network congestion

❖ sender limits transmission:



sender sequence number space

cwnd

last byte ACKed

sent, but not-yet ACKed ("in-flight")

last byte sent

available but not used

`LastByteSent-LastByteAcked ≤ Min(Rwnd, cwnd)`

❖ *TCP sending rate:* cwnd bytes every RTT

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

**On average, over a RTT time period.**

❖ *Approach: Feedback control loop:* sender increases transmission rate (window size), probing for usable bandwidth, until loss occurs.

▪ *additive increase:* increase `cwnd` by **1** MSS every RTT until loss is detected.

▪ *multiplicative decrease:* cut `cwnd` in half after loss (perceived congestion).

AIMD saw tooth behavior: probing for bandwidth

additively increase window size …

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

`cwnd`: TCP sender congestion window size

time

• Why AIMD?
  • AI: Gradual capacity/bandwidth probing with low potential for loss; improved fairness.
  • MD: Since increased congestion is more catastrophic, reduce it more aggressively.

**Q1:** What should a congestion window (cwnd) start with?

Initially we can set **cwnd** = 1 MSS (minimum/safe value)

```
CWND = (MSS > 2190)? (2*MSS): (MSS > 1095)? (3*MSS): (4*MSS)
```
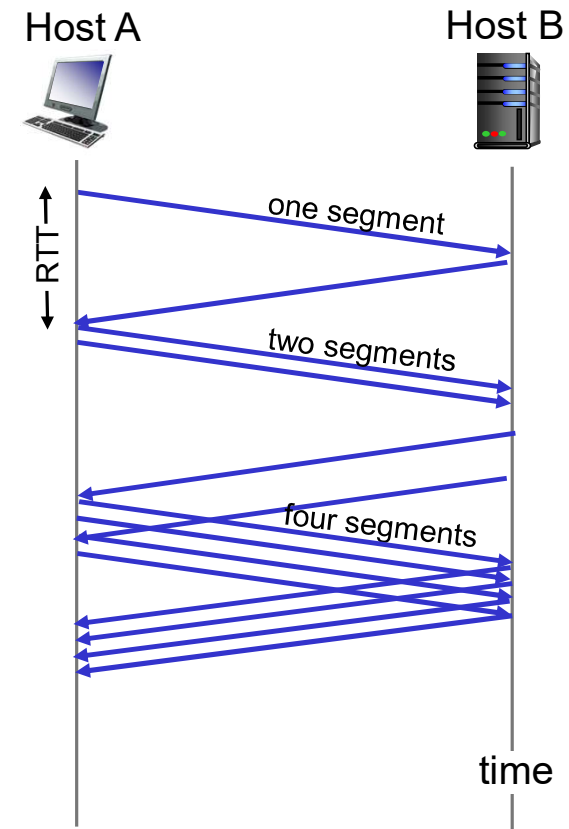
**Q2:** How exactly to grow/shrink the CWND?

❖When connection begins, increase the rate exponentially /linearly until a ***certain threshold*** or the ***first loss event***:

- Exponentially: double **cwnd** every RTT
- Linearly: increment **cwnd** by 1/cwnd for every ACK received.

```
If cwnd < ssthresh then
        Each time an Ack is received:
        cwnd = cwnd + 1
else: /* cwnd >= ssthresh */
        Each time an Ack is received:
        cwnd = cwnd + 1 / [ cwnd ]
endif
```

Host A          Host B

RTT

one segment

two segments

four segments

time

*summary:* **initial rate is small (slow) but ramps up exponentially fast**

**Q3:** When should the exponential increase switch to linear increase?

❖ When we expect to be reaching near to the network capacity.

# Implementation:

❖ Maintain variable `ssthresh`

❖ on loss event, `ssthresh` is set to 1/2 of `cwnd`

`ssthresh = max(cwnd/2, 2*MSS)`

```
If Timeout then: //or 3D-ACK (Tahoe) then:
//Indicates packet loss: (MD)
        ssthresh = max(cwnd /2, 2*MSS)
        cwnd = 1 or cwnd = ssthresh + 3
Else: //regular ACK (@every ack)
        If cwnd < ssthresh then:
                cwnd = cwnd + 1
        else: /* cwnd >= ssthresh */
                cwnd = cwnd + 1 / [ cwnd ]
        endif
endif
```
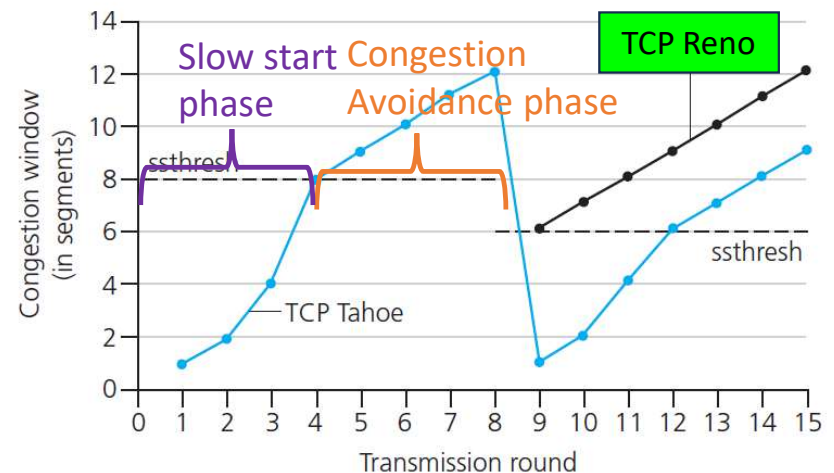


Slow start phase   Congestion Avoidance phase   TCP Reno
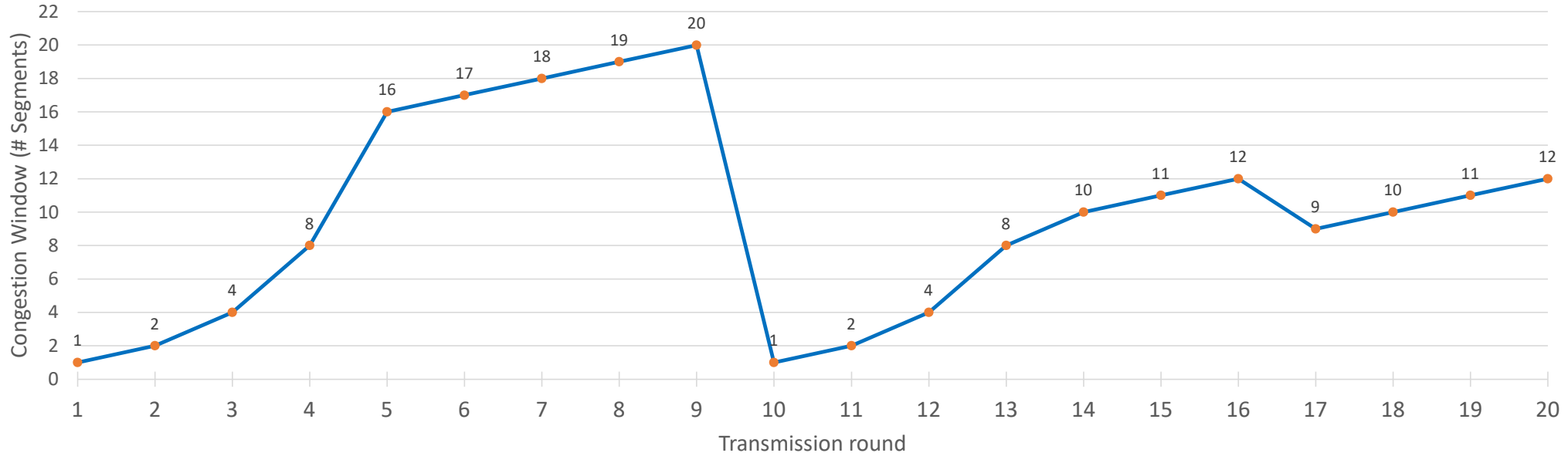
*Why should we modify the **ssthresh** on loss?*

## Consider: RTT=1s; MSS=1000 bytes



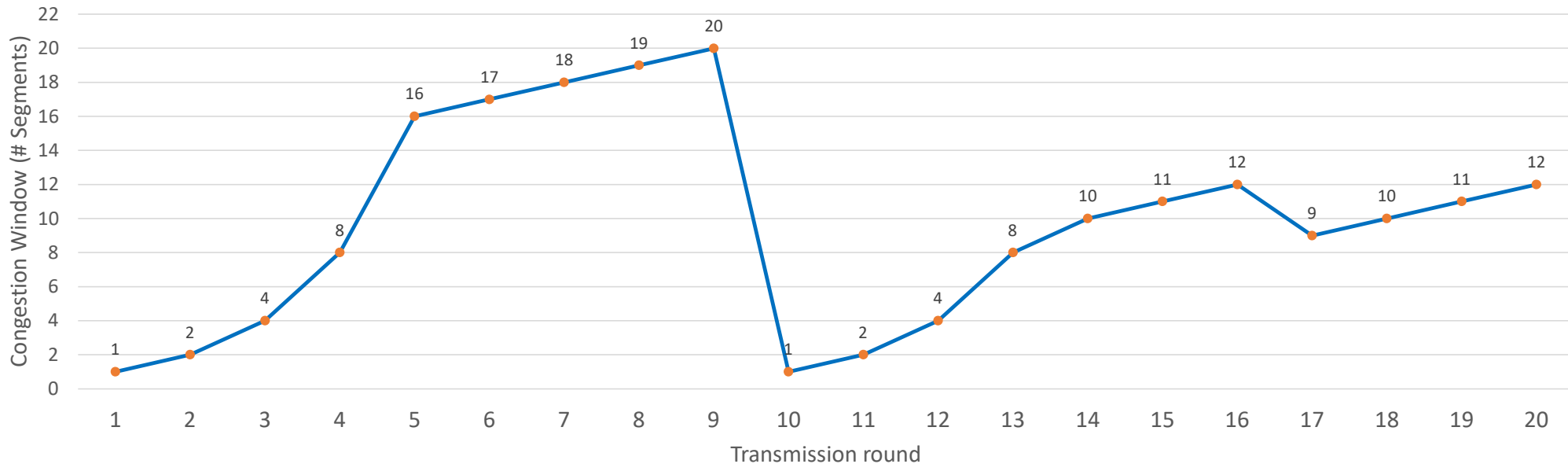**Q1**: What is the initial ssthresh value if any?

**Q2**: How many congestion events occurred?

**Q3**: Identify the slowstart and CA phases?

**Q4**: ssthresh value at Transmission round 12 and 18 respectively?

## Consider: RTT=1s; MSS=1000 bytes



**Q5**: Name the congestion control mode used?

**Q6**: How many segments were delivered until first congestion event?

**Q7**: Average throughput in first slow-start phase?

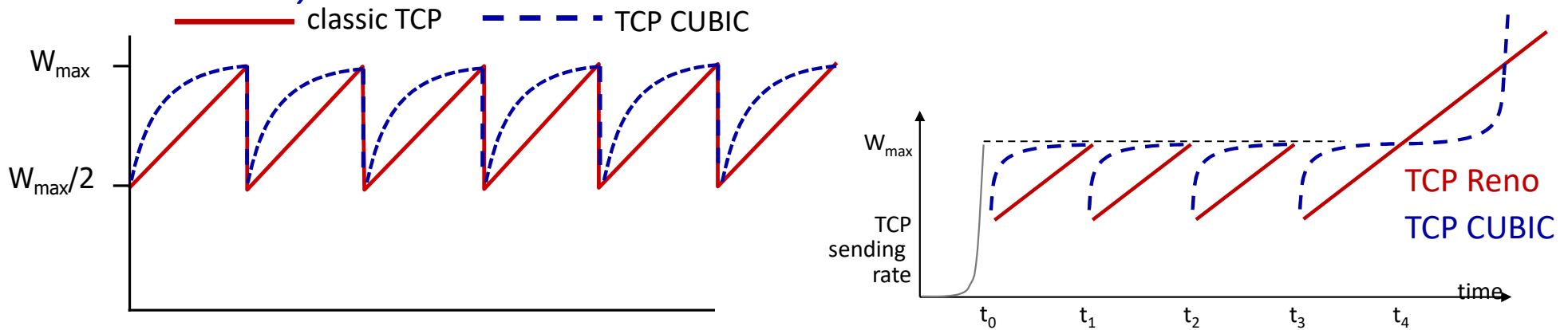**Q8**: Average throughput in first congestion avoidance phase?

# TCP Variants: TCP Tahoe to TCP BBR

- **TCP Tahoe:** Introduced the congestion control mechanisms: Slowstart, Congestion Avoidance.
- **TCP Reno**: Fast retransmit and fast recovery are added in addition to previous congestion control mechanisms.
  - Has other features -- header compression (if ACKs are being received regularly, omit some fields of TCP header).
  - Delayed ACKs -- ACK only every other segment.
- **TCP NewReno**: Modified and Advanced Fast recovery mode.
  - Can handle multiple packet losses -- Extension likewise to TCP SACK.
- **TCP Vegas**: Rate/Delay based Congestion Detection.
  - New Congestion Avoidance and Retransmission mechanism.
- **TCP CUBIC**: Provides Faster Congestion response in high speed networks.
  - Default TCP congestion control setting in Linux.
- **TCP BBR**: Recent Extension by Google (Also rate/delay based).
  - Supported in Linux v 4.9 onwards.
- Several other variants like **TCP Fast, Hybla, Illinois, Veno, WestWood, etc.** are not discussed.
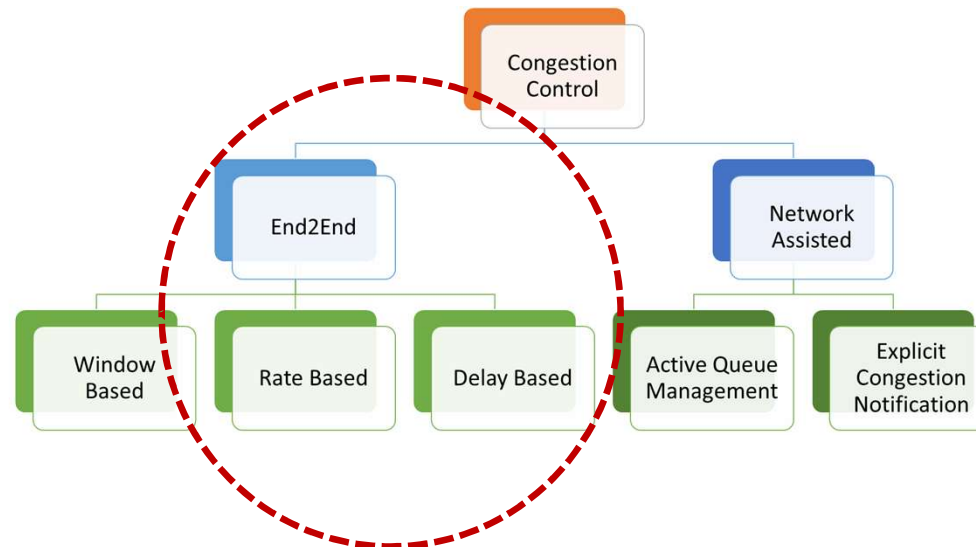
https://en.wikipedia.org/wiki/TCP_congestion_control

- **Is there a better way than AIMD to "probe" for usable bandwidth?**
- **Insight/intuition:**
  - $W_{max}$: sending rate at which congestion loss was detected.
  - congestion state of bottleneck link **probably (?)** hasn't changed much.
  - after cutting window in half on loss, initially ramp to to $W_{max}$ *faster*, but then approach $W_{max}$ more *slowly.*
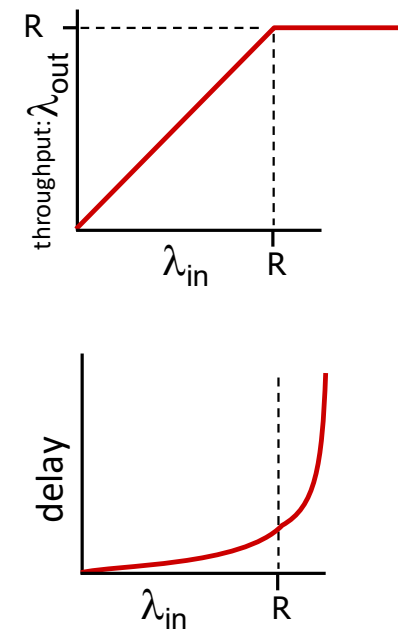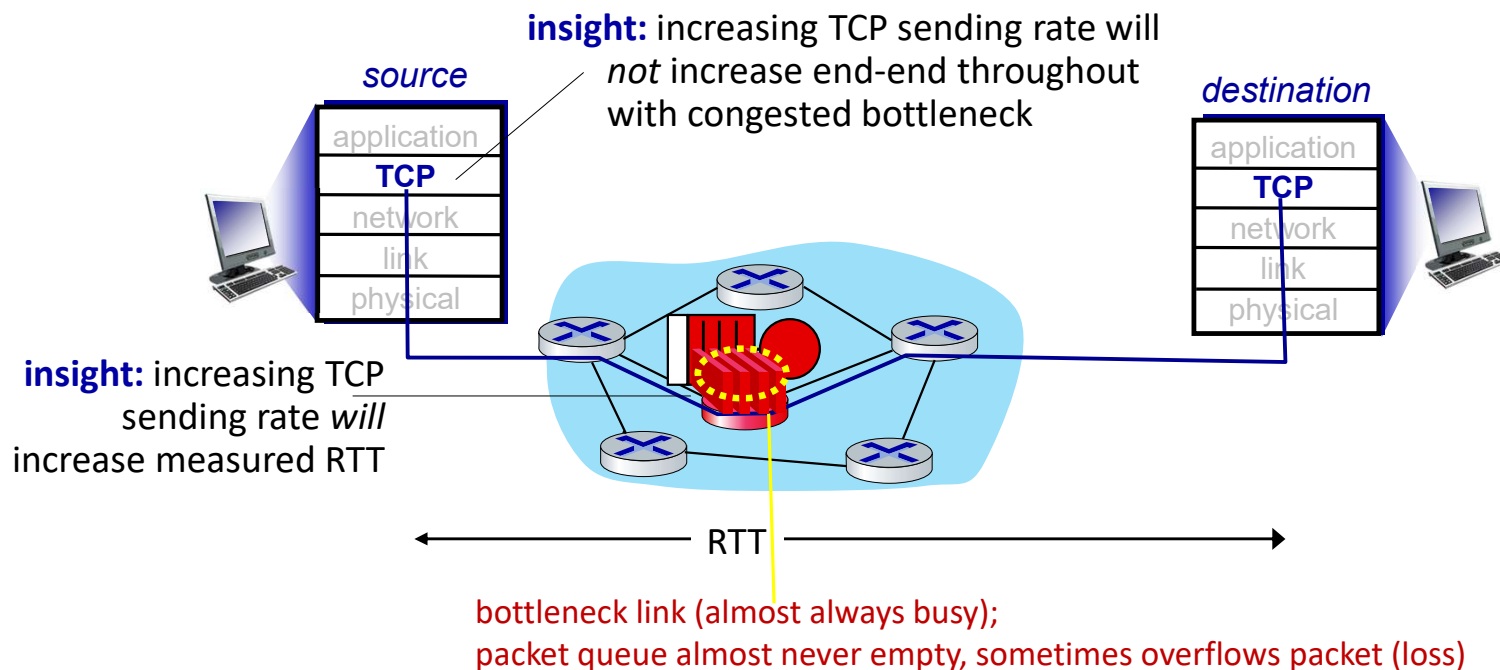


- **K – a tunable point in time when TCP window size will reach $W_{max}$.**
- **increase W as a function of the *cube* of the distance between current time & K**
  - larger increases when further away from K
  - smaller increases (cautious) when nearer K

# Delay/Rate based Congestion Control
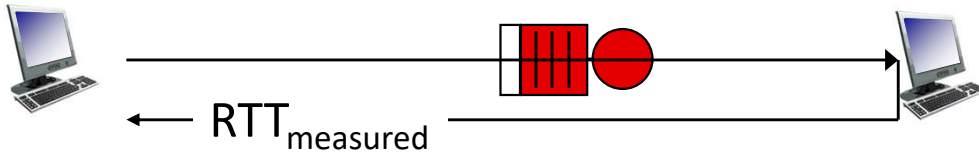
# TCP AND THE CONGESTED "BOTTLENECK LINK"

- understanding congestion: useful to focus on congested bottleneck link

**insight:** increasing TCP sending rate will *not* increase end-end throughout with congested bottleneck

*source*

application
**TCP**
network
link
physical

*destination*

application
**TCP**
network
link
physical

**insight:** increasing TCP sending rate *will* increase measured RTT

RTT

bottleneck link (almost always busy);
packet queue almost never empty, sometimes overflows packet (loss)

*Goal:* *"keep the end-end pipe just full, but not fuller"*

throughput: $\lambda_{out}$

R

$\lambda_{in}$    R

delay

$\lambda_{in}$    R

Keeping sender-to-receiver pipe "just full enough, but no fuller": keep bottleneck link busy transmitting, but avoid high delays/buffering

$\leftarrow$ RTT$_{measured}$

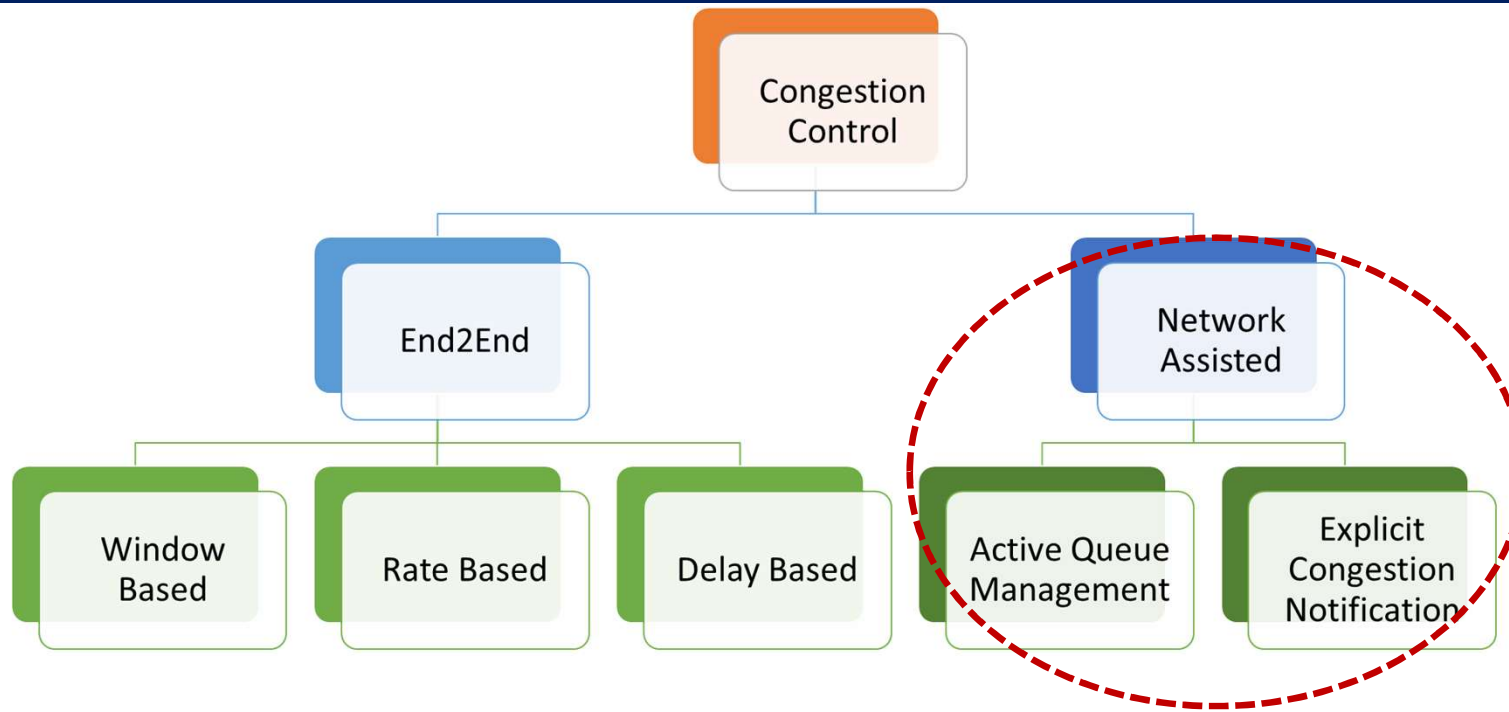$$\text{measured throughput} = \frac{\text{\# bytes sent in last RTT interval}}{\text{RTT}_{measured}}$$

## Delay-based approach:

- congestion control without inducing/forcing loss $\rightarrow$ Maximizes throughput

- a number of deployed TCPs take a delay-based approach

  - BBR deployed on Google's (internal) backbone network.
  - RTT$_{min}$ : minimum observed RTT (i.e. RTT for uncongested path)
  - uncongested throughput with congestion window `cwnd` is cwnd/RTT$_{min}$

IF *measured throughput* "very close" to  uncongested throughput:
    increase `cwnd` linearly          /* since path not congested */
ELSE IF ***measured throughput*** "far below" uncongested throughout:
    decrease `cwnd`  linearly         /* since path is congested */

- Split the responsibility of congestion control between end hosts and routers.

1. Router monitors congestion & explicitly notifies end-hosts when congestion is about to occur.

2. In response, the source adjusts transmit rate upon such notifications.

# Random Early Drop (RED)

- Each router monitors its queue length.

- IF the *queue length* > **threshold**, THEN:

  - Drop packets (implicit notification) to indicate network congestion.

  Or

  - Explicitly notify source by marking "Explicit Congestion Notification" bit.

- Note that packets are dropped much earlier than usual -- before buffer resources are exhausted completely -- *so drops are fewer*

  - Bursty drops are also reduced.

# DETAILS OF RED

- The principle is to drop the packet with some "**drop probability**" when the queue length exceeds a certain "***drop level***".

**Q1:** Should we use the Sample Queue Length or the average queue length?

- Instead of a sample queue length, average queue length (more accurately captures notion of congestion) is considered.

$$AvgLen = (1-weight)*AvgLen + weight * SampleLen$$
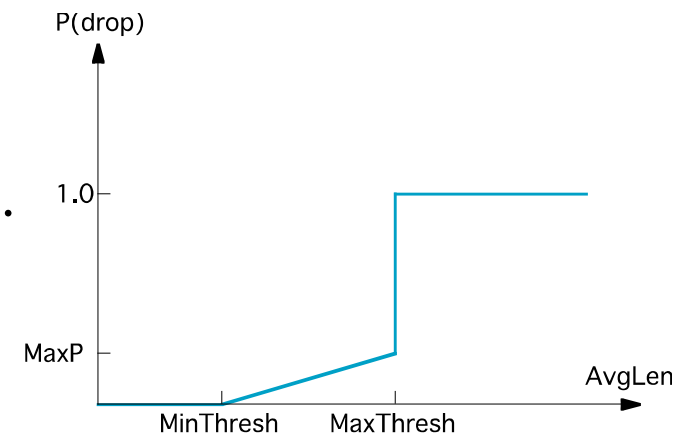
**Q2:** How to effect the packet drops?

- Two thresholds : **MinThreshold** and **MaxThreshold**.

    If AvgLen <= MinThreshold queue the packet

    If AvgLen >= MaxThreshold drop the arriving packet.

    If MinThreshold <= AvgLen <= MaxThreshold, then,

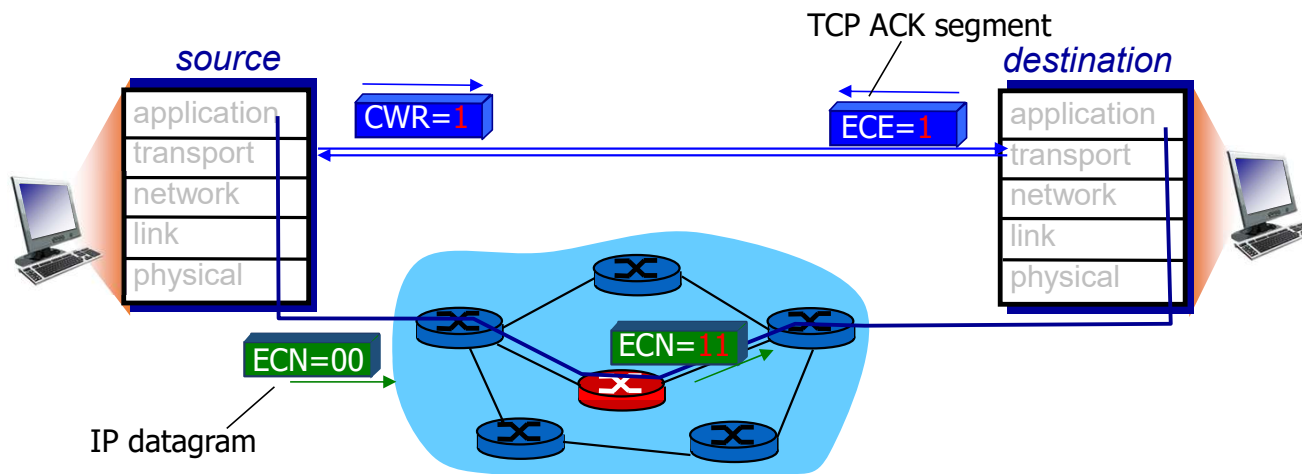    calculate a ***drop probabilty P*** and drop the arriving packet with the probability P.

P(drop)

1.0

MaxP

MinThresh     MaxThresh

AvgLen

Reading more on RED: http://dl.acm.org/citation.cfm?id=169935

**Q3:** Can we do better than just dropping the packets?

- Instead of packet drop, use the same packets to carry and notify the indication of congestion.

- two bits in IP header (ToS field) marked *by network router* to indicate congestion.

- receiver (seeing congestion indication in IP datagram) ) sets ECE bit on receiver-to-sender ACK segment to notify sender of congestion.
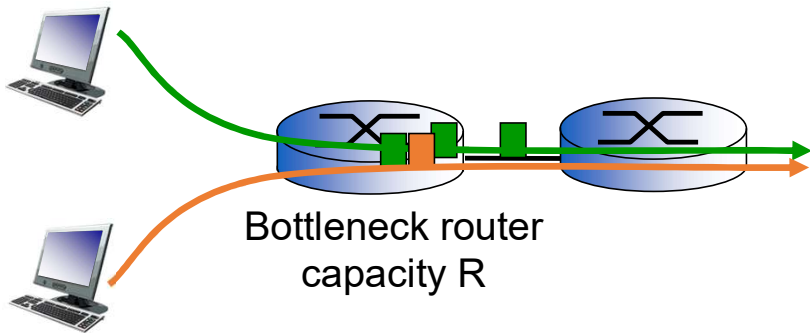


Reading more on ECN: https://dl.acm.org/doi/pdf/10.17487/RFC2481
https://www.cse.wustl.edu/~jain/papers/ftp/cr1.pdf

# TCP Fairness

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

TCP connection 1
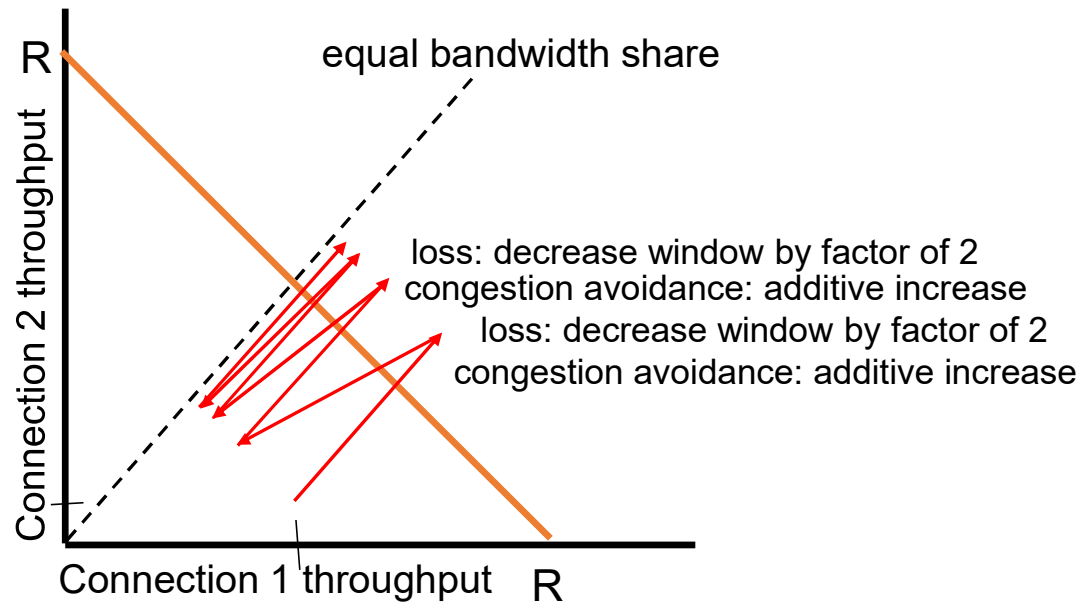
Bottleneck router
capacity R

TCP connection 2

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

Connection 1 throughput   R

R

## Fairness and UDP

❖ multimedia apps often do not use TCP

  ▪ do not want rate throttled by congestion control

❖ instead use UDP:

  ▪ send audio/video at constant rate, tolerate packet loss

## Fairness, parallel TCP connections

❖ application can open multiple parallel connections between two hosts

❖ web browsers do this

❖ e.g., link of rate R with 9 existing connections:

  ▪ new app asks for 1 TCP, gets rate R/10

  ▪ new app asks for 11 TCPs, gets R/2

*A Fundamental Advancement in TCP!*
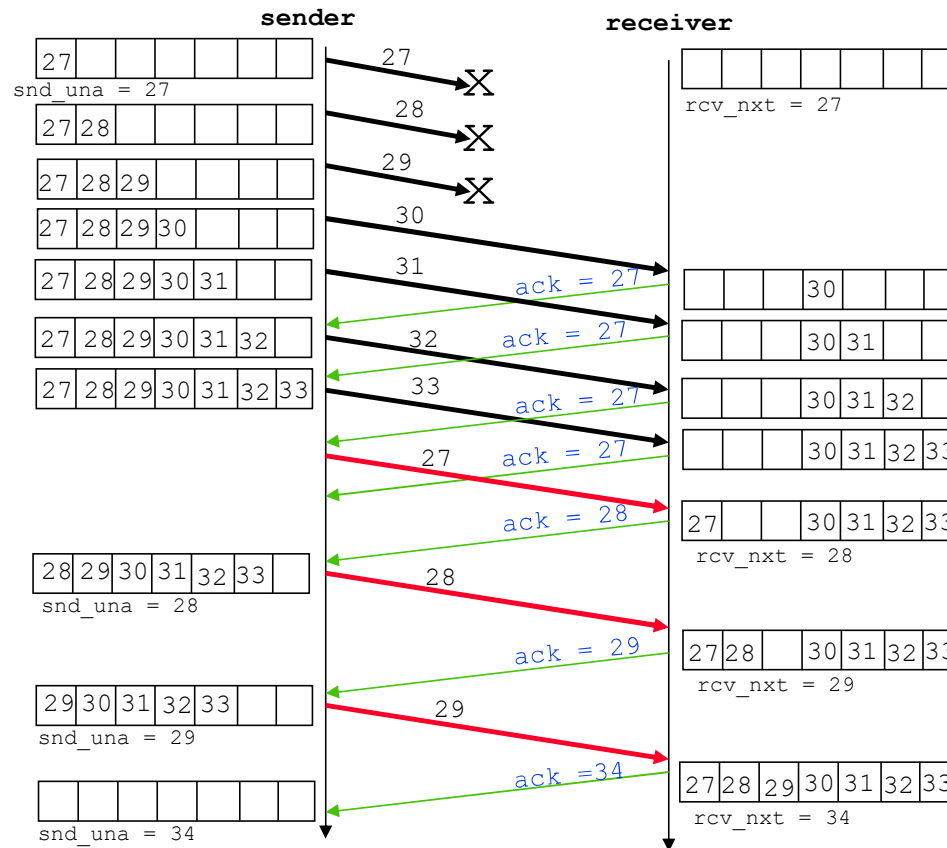*<Not covered in the text books!>*


*TCP SACK + Selective Repeat Retransmission Policy.*


https://tools.ietf.org/html/rfc2018

... 

# SELECTIVE ACKNOWLEDGEMENT (RFC 2018)

- SACK = <u>S</u>elective <u>Ack</u>nowledgement ('New' TCP option)

- Reasons for SACK
  - TCP may perform poorly when multiple packets are lost
  - Data may be retransmitted unnecessarily by go-back-N (TCP Tahoe)
  - Or one missing segment retransmitted per round trip

- Receiver's part of SACK
  - Acknowledges non-contiguous blocks of data with SACK blocks
  - SACK blocks: sequence number left/right edges of isolated data
  - May drop data that is "SACKed"

- Sender's part of SACK
  - SACKed blocks marked in retransmission buffer
  - Marked segments not retransmitted, "holes" retransmitted
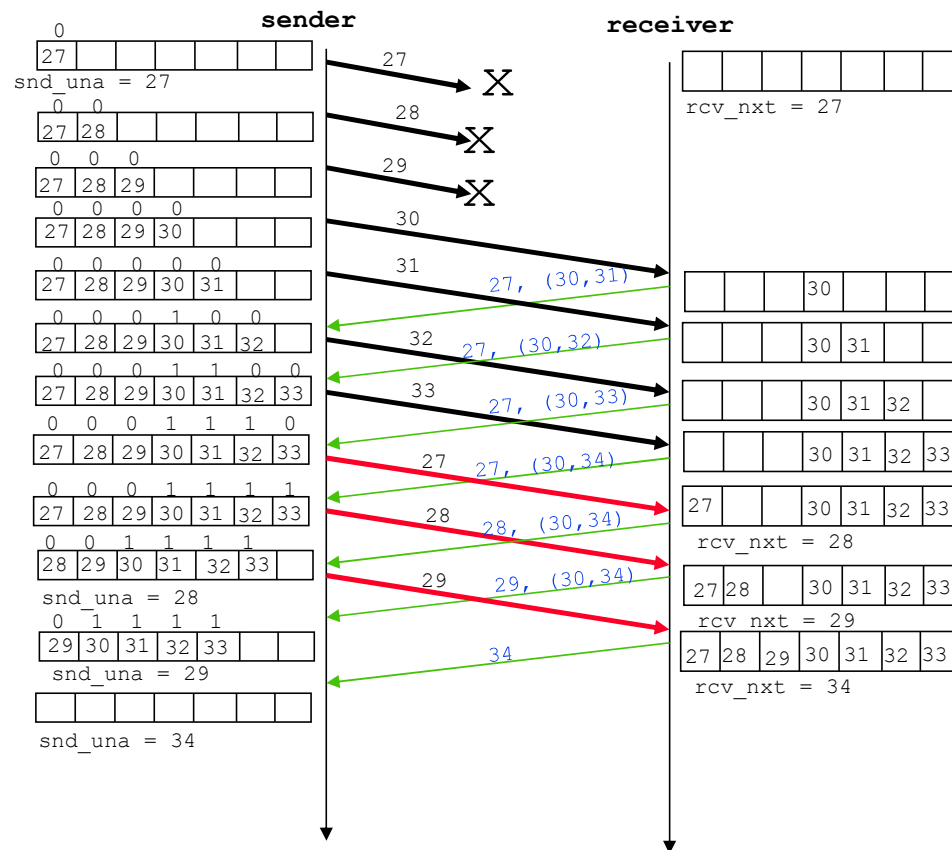  - Blocks unmarked after retransmission timeout

- Window size = 7
- Multiple packets lost in one RTT

- retransmit pkt 27 on receiving 2(3 total) dup acks.
- Pkt. 28 can be retransmitted only after receiving ack of pkt 27
  - 1 RTT wait

- receiver waits longer to deliver packets to application in sequence

- Still retransmit pkt 27 on receiving 2 (3 total) dup (s)acks.

- Pkt. 28 can be retransmitted right away, knowing that 3 packets have been received subsequently
  - only wait an inter-packet transmit time,

- ε reduction in latency

- For High Bandwidth-Delay Connections
  - Accurate Round-Trip Time estimation
  - Overcome limit on window size
  - Dealing with Packet loss

- Use TCP Options
  - Timestamp Extension
  - Protecting from Sequence number wraparound (32 bits ~36 Gbits)
  - Extensions for Larger Windows (Window Scaling)