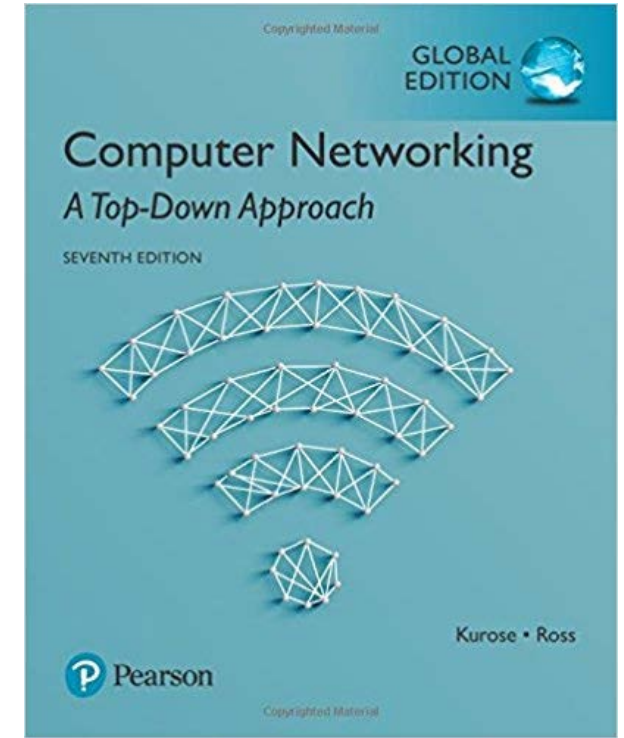# Chapter 3
# Transport Layer part 4

School of Computing
Gachon Univ.
Joohyung Lee

Most of slides from J.F Kurose and K.W. Ross. And, some slides from Prof. Joon Yoo
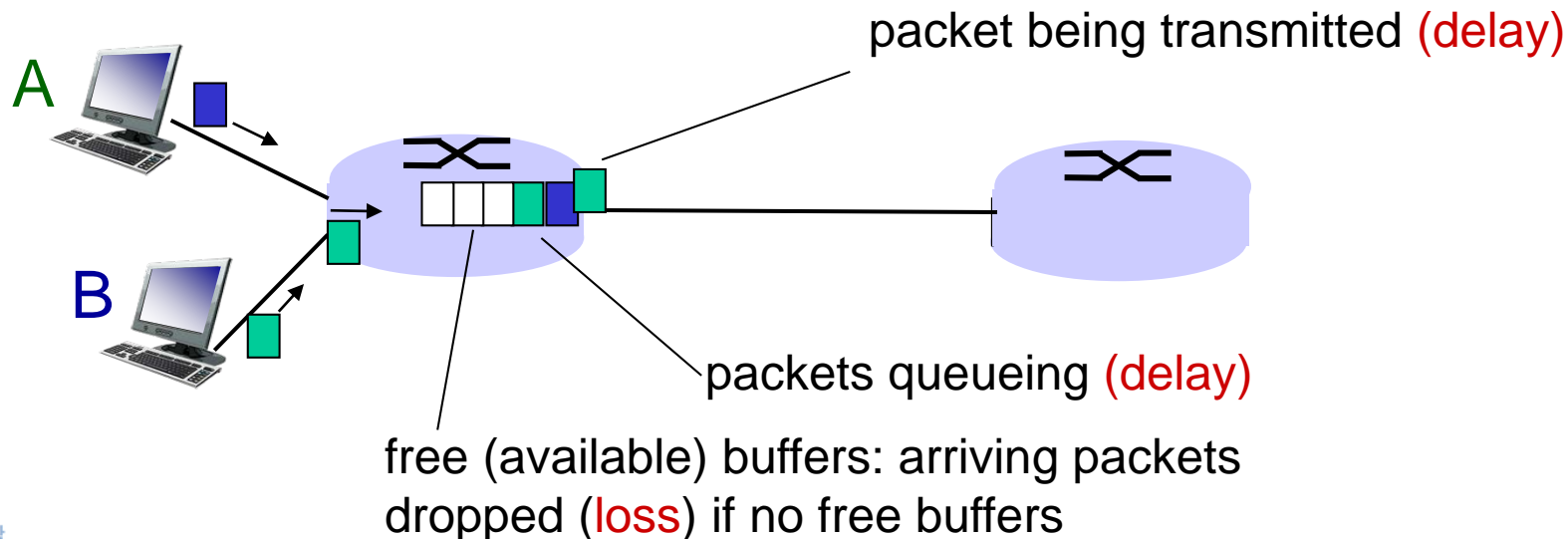
*Computer Networking: A Top Down Approach*

# Chapter 3 outline

# How do loss and delay occur?

## packets *queue* in router buffers

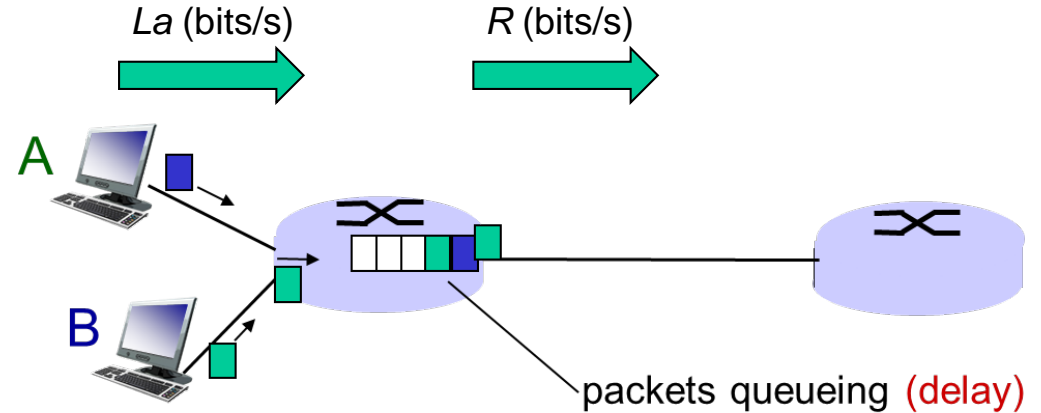❖ packet arrival rate to link (temporarily) exceeds output link capacity

❖ packets queue, wait for turn

❖ Unlike other delays ($d_{proc}$, $d_{trans}$, and $d_{prop}$), queueing delay ($d_{queue}$) can vary from packet to packet

**queueing delay : variation** .
( )

A

packet being transmitted (delay)

B

packets queueing (delay)

free (available) buffers: arriving packets
dropped (loss) if no free buffers

# Queueing delay (Ch. 1.4.2)

- ❖ **L:** packet length (bits)
- ❖ **a**: average packet arrival rate (packets/s)
- ❖ **R**: link bandwidth (bps)

*La* (bits/s)   *R* (bits/s)

A

B

packets queueing (delay)

link capacity        queueing delay

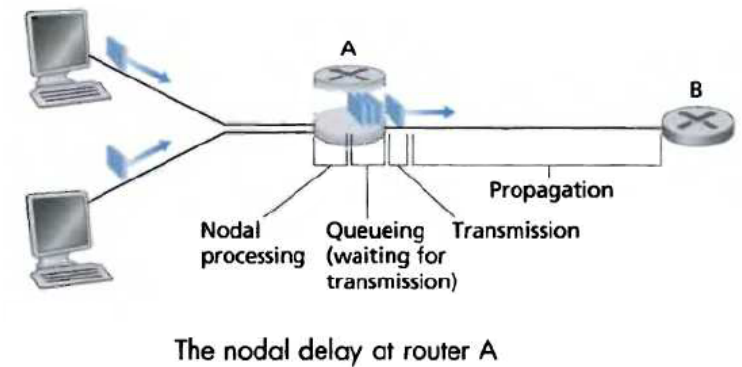- ❖ **La/R = $\rho$ : Traffic Intensity**
- ❖ $\rho > 1$: more "work" arriving than can be serviced; queue will tend to increase without bound; average delay infinite.
  - ▪ This should never happen!
- ❖ Design the system so that $\rho \leq 1$

queue

# Principles of congestion control

*congestion*:

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

❖ Symptoms (signs):
  = **queueing delay**
  - long delays (queueing in router buffers)
  - **lost packets** (buffer overflow at routers)

❖ **Congestion Control:**
  - ways to treat the cause of **network congestion**
  - different from **flow control**!
    - flow control considered the *receiver*'s capacity



The nodal delay at router A

# Causes/costs of congestion: scenario 1

- ❖ two senders, two receivers
- ❖ one router, *infinite buffers*
  - ▪ packet loss?
- ❖ output link capacity: **R**
  - ▪ each flow shares R/2
- ❖ no retransmission

original data: $\lambda_{in}$

Host A

throughput: $\lambda_{out}$

30

unlimited shared
output link buffers

R

Host B

30

50



R/2
queue

- ❖ maximum per-connection throughput: R/2

- ❖ **large queuing delays** as arrival rate, $\lambda_{in}$, approaches capacity

미래창조과학부 SW중심대학
가천대학교 소프트웨어학과
Gachon University Department of Software

# Causes/costs of congestion: scenario 2

- ❖ one router, *finite* buffers
  - packet loss?
- ❖ sender retransmission of timed-out packet
  - application-layer input = application-layer output: $\lambda_{in} = \lambda_{out}$
  - transport-layer input includes *retransmissions* : $\lambda'_{in} \geq \lambda_{in}$
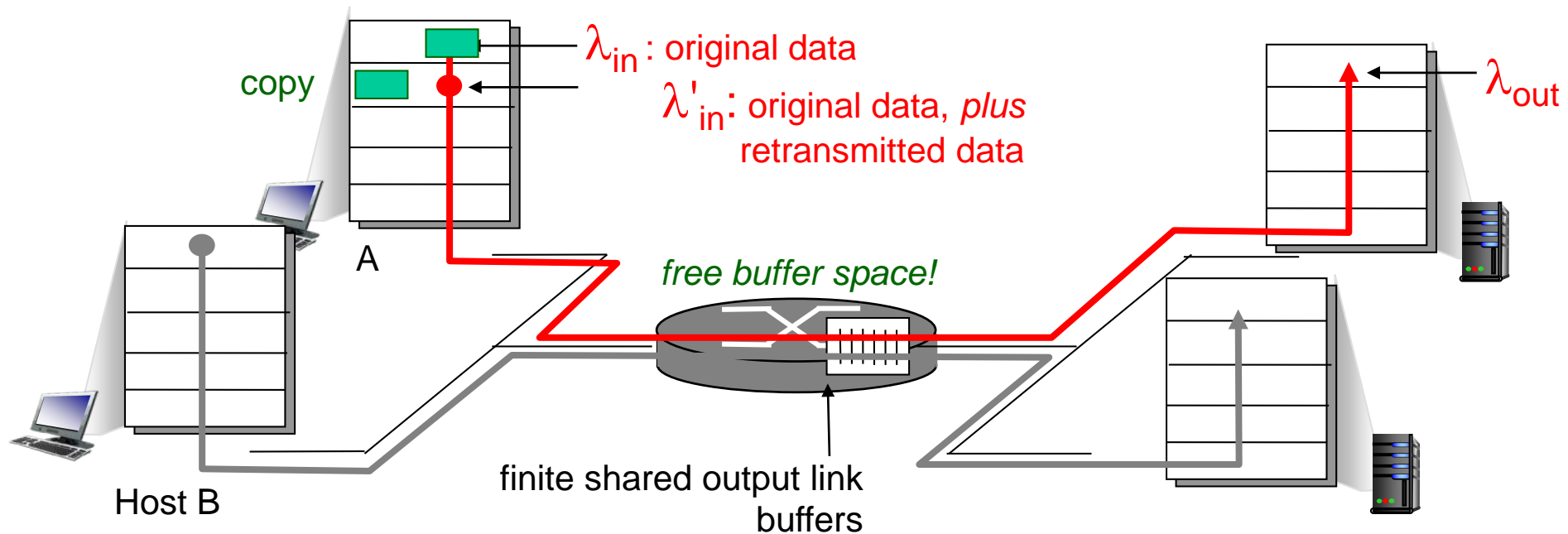
$\lambda_{in}$ : original data

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

$\lambda_{out}$

Host A

Host B

finite shared output link buffers

미래창조과학부 SW중심대학
가천대학교 소프트웨어학과
Gachon University Department of Software

# Causes/costs of congestion: scenario 2
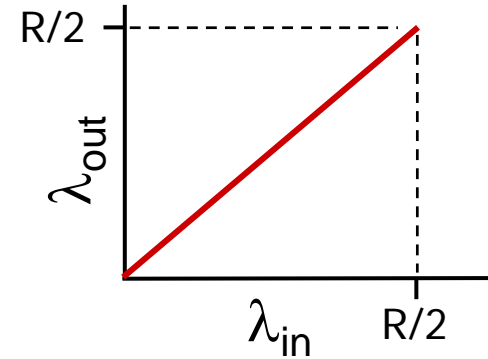
idealization1: perfect knowledge

❖ sender sends only when router buffers available

- $\lambda'_{in} = \lambda_{in} = \lambda_{out}$

loss



$\lambda_{in}$ : original data

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

$\lambda_{out}$

copy

A

free buffer space!

finite shared output link buffers

Host B

미래창조과학부 SW중심대학
가천대학교 소프트웨어학과
Gachon University Department of Software

# Causes/costs of congestion: scenario 2

*Idealization2: known loss*
packets can be lost, dropped at router due to full buffers

❖ sender only resends if packet *known* to be lost

- $\lambda'_{in} \geq \lambda_{in}$

copy

A

Host B

$\lambda_{in}$ : original data

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

*no buffer space!*

$\lambda_{out}$
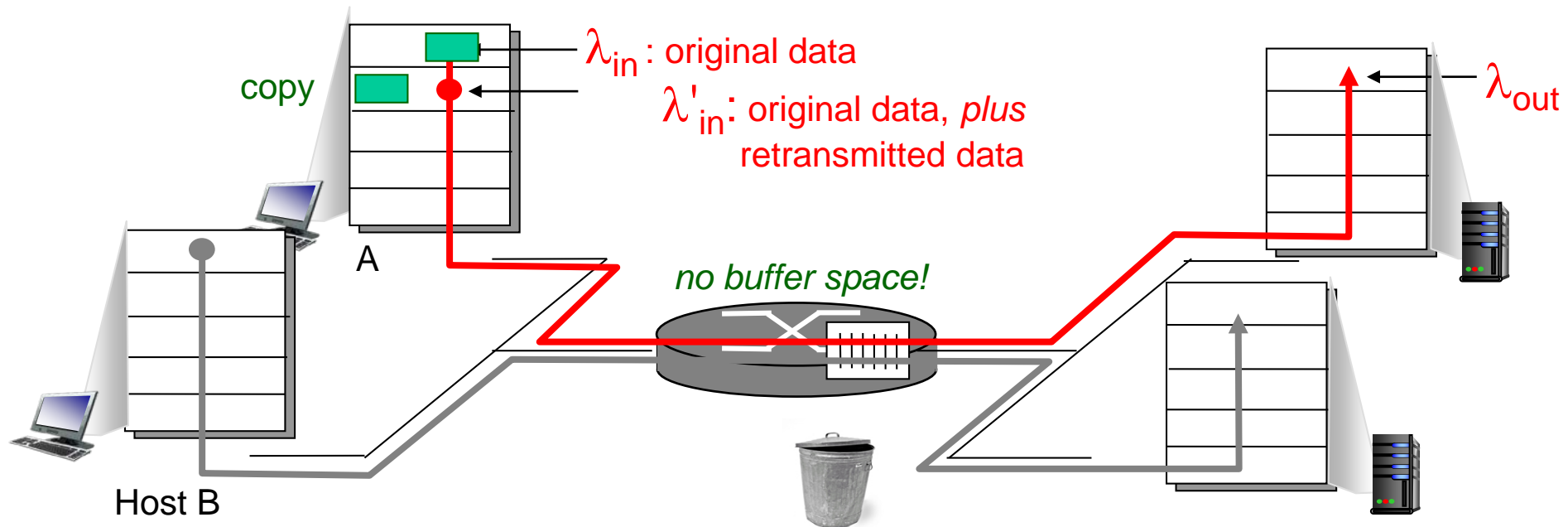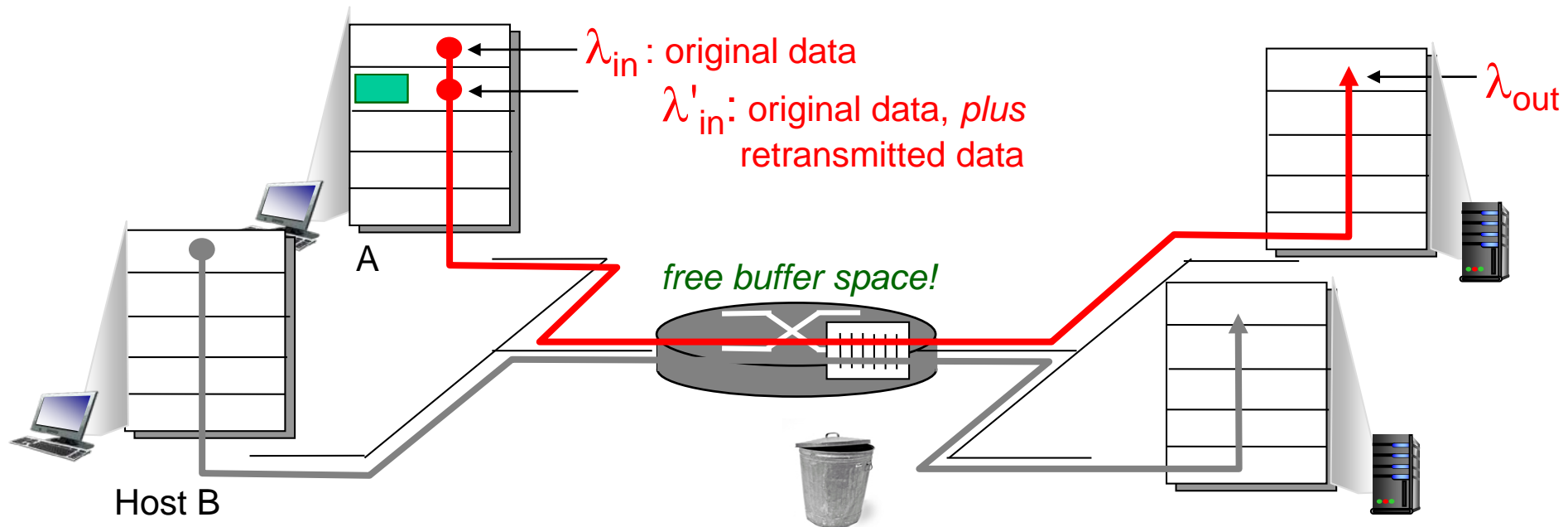
# Causes/costs of congestion: scenario 2

*Idealization2: known loss*

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

- ❖ sender only resends if packet *known* to be lost
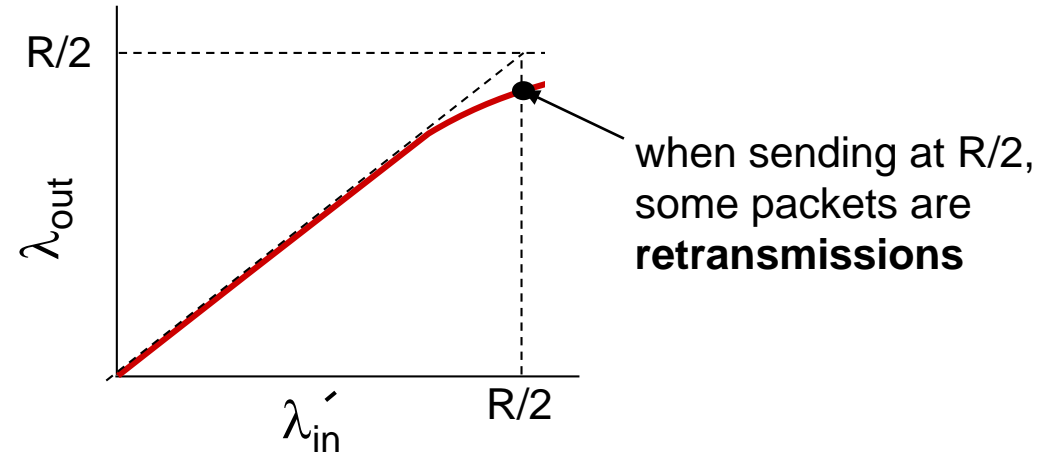  - ■ $\lambda'_{in} \geq \lambda_{in}$

when sending at R/2, some packets are **retransmissions**

$\lambda_{in}$ : original data

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

$\lambda_{out}$

A

*free buffer space!*

Host B

# Causes/costs of congestion: scenario 2

## *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



when sending at R/2, some packets are retransmissions including **duplicated** that are delivered!

timeout

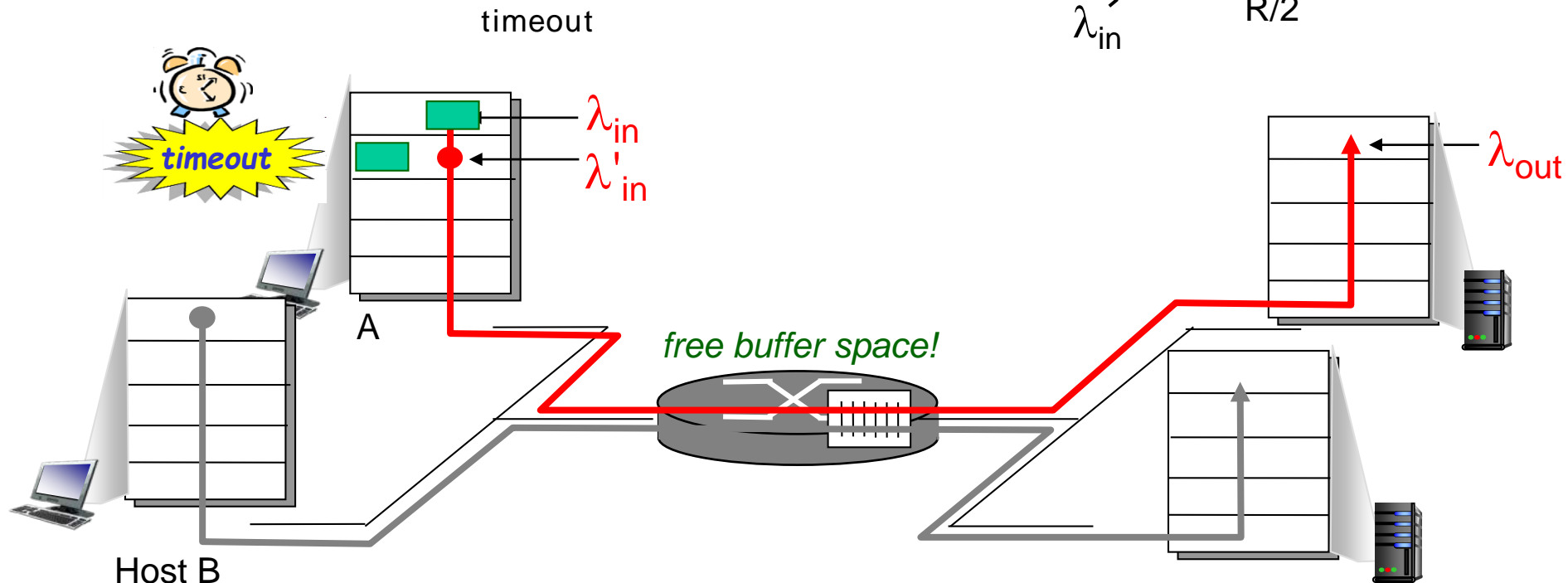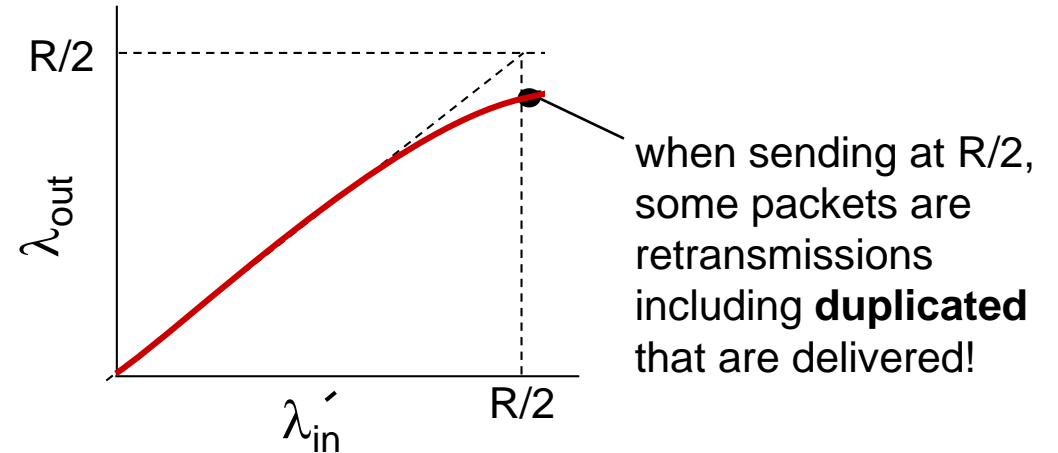$\lambda_{in}$

$\lambda'_{in}$

timeout

A

free buffer space!

$\lambda_{out}$

Host B

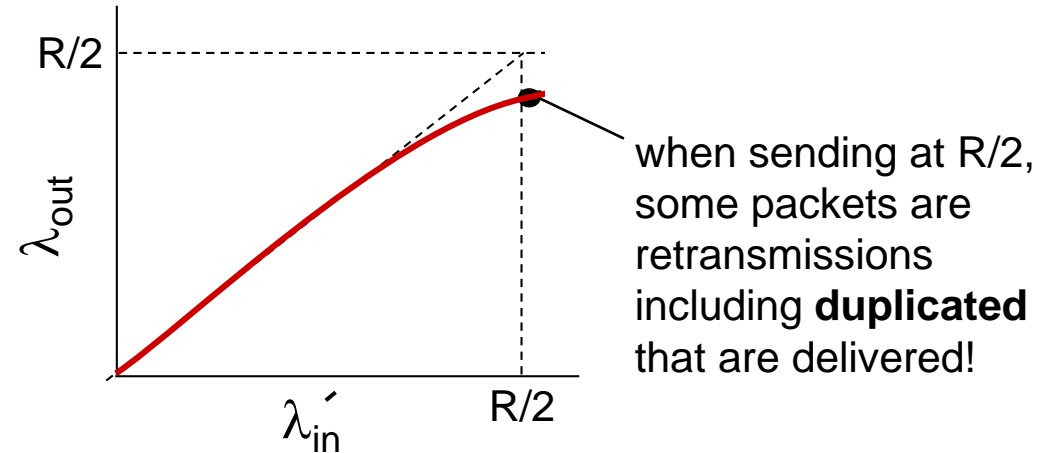# Causes/costs of congestion: scenario 2

## *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



when sending at R/2, some packets are retransmissions including **duplicated** that are delivered!

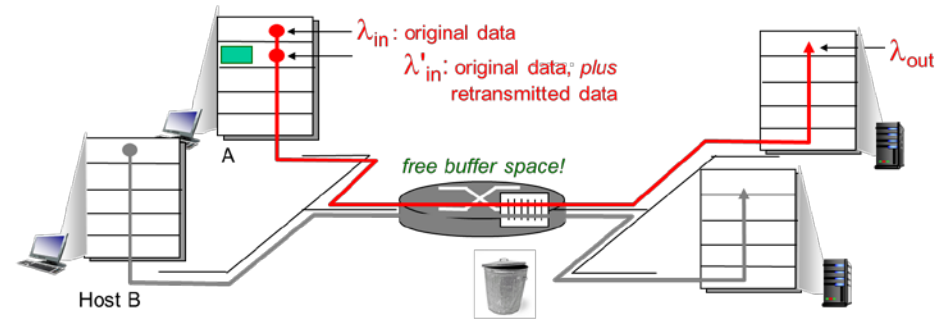## "costs" of congestion:

❖ large queueing delays

❖ more work (retrans) for given "goodput"

❖ unneeded retransmissions: link carries multiple copies of pkt
  ▪ decreasing goodput

미래창조과학부 SW중심대학
가천대학교 소프트웨어학과
Gachon University Department of Software

# From single-hop to multi-hop path

❖ Until now (scenario 2), we only considered a **single-hop** path congestion



❖ But, generally in the Internet, the network path is **multi-hop** (Scenario 3)

❖ Congestion can happen at any link along the multi-hop path



congestion          goodput

# Causes/costs of congestion: scenario 3



another "cost" of congestion:

❖ when packet dropped, any transmission capacity used along the path for that packet was wasted!
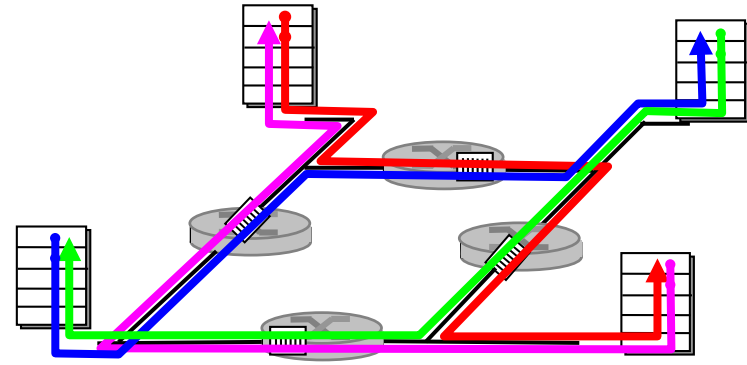
# Chapter 3 outline

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

:        **congestion**

(    )

Transport Layer     3-15

# TCP congestion control

- ❖ TCP/IP protocol stack was operational in early 1980s                    **link**
- ❖ But, Internet was suffering from congestion collapse          **capacity**
  - hosts would send their packets into the Internet as fast as received window would allow (only flow control)
  - congestion would occur at some router and packet drop happens, and hosts would retransmit their packets, resulting in even more congestion
- ❖ TCP congestion control introduced by Van Jacobson in later 1980s
  - Each host (end systems) determines how much capacity is available in the network



**sending rate**

# TCP congestion control

router                                                    client

❖ End-to-end congestion control   link capacity        traffic intensity

❖ Sender needs to determine **network capacity**, which changes over time
                                         n              (n=1    1    ack    )

❖ Sender limits the **sending rate** as a function of available network capacity (i.e., network congestion)

  ▪ e.g., little congestion – increase **send rate**

  ▪ e.g., congestion along path – decrease **send rate**

❖ The only feedback to sender: **ACK**        ack        congestion

  ▪ ACK arrival: a packet has arrived safely at receiver

  ▪ ACK timeout: has not arrived   -> congestion

# TCP congestion control

❖ Questions

1. **How does the sender limit the rate?** <sup>n</sup>

2. How does the sender know there is congestion on path between itself and destination?

3. What algorithm should sender use to change send rate as function of congestions?

# TCP Congestion Control: congestion window

n(segment ) TCP
cwnd(segment or byte -> segment * byte )
*sender sequence number space*



last byte ACKed

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

last byte sent

❖ sender limits transmission:
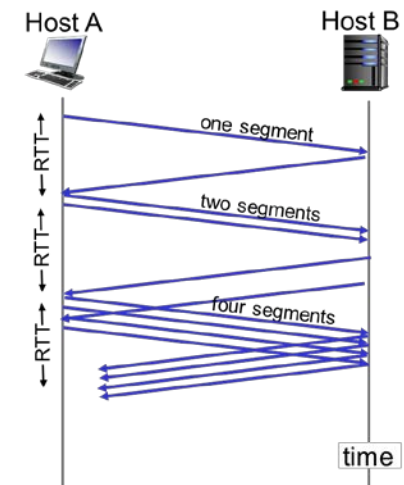
$$LastByteSent - LastByteAcked \leq cwnd$$

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

*TCP sending rate:*

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

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



Host A          Host B

one segment

two segments

four segments

time

# TCP Congestion Control: congestion window

***Congestion control(cwnd)*** *and* ***Flow control (advertised wnd)*** *determine TCP sending rate*

**window**

:

Send Window = MIN(flow control window, **congestion window**)

Note: advertised window is not so dynamic!

**= advertised window**

미래창조과학부 SW중심대학
가천대학교 소프트웨어학과
Gachon University Department of Software

# TCP congestion control

❖ Questions

1. How does the sender limit the rate?

2. **How does the sender know there is congestion on path between itself and destination?**

3. What algorithm should sender use to change send rate as function of congestions?

# Network Congestion

congestion

❖ How can the TCP sender know there is congestion?
  ▪ Where does congestion actually happen?
  ▪ However, network layer (e.g., router) does not send explicit message of congestion

❖ Loss event!
  ▪ When there is excessive congestion, one (or more) router buffers along the path overflows, causing packet drop – **loss**!
  ▪ Indication of loss
    • Timeout
    • Receipt of three duplicate ACKs
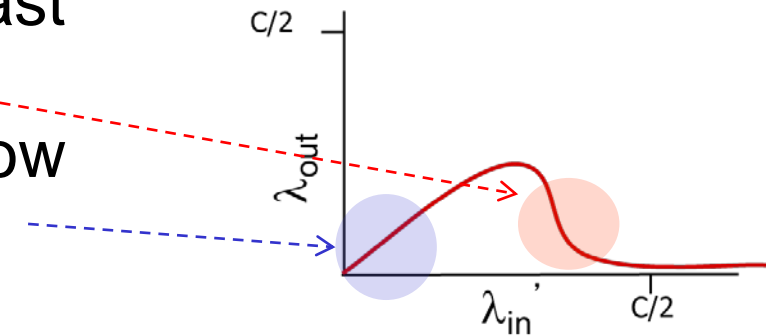
# TCP congestion control

❖ Questions

1. How does the sender limit the rate?

2. How does the sender know there is congestion on path between itself and destination?

3. **What algorithm should sender use to change send rate as function of congestions?**

미래창조과학부 SW중심대학
가천대학교 소프트웨어학과
Gachon University Department of Software

# TCP congestion control principles

❖ **How should TCP sender determine send rate (or cwnd)?**

- If TCP senders collectively send too fast
  – network congestion may happen
- If TCP senders cautiously send too slow
  – underutilize the network
  utilization



❖ **TCP principles**

- A lost segment implies congestions, and hence, the TCP sender's rate should be decreased
- An acknowledged segment indicates successful delivery, and hence, the sender's rate can be increases when ACK arrives
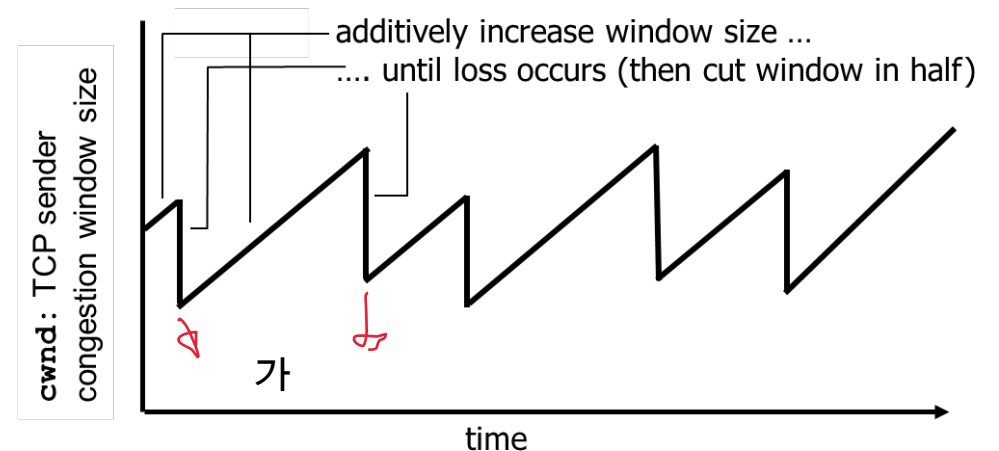
# TCP congestion control principles

❖ **Bandwidth probing**
  - TCP sender increases its rate in response to arriving ACKs until loss event occurs – then decrease rate
  - TCP sender increases its rate to probe for rate at which congestion begins
  - Then backs off from that rate, and then begins probing again to see if the congestion rate has changed.

❖ **TCP Congestion Control Algorithm**
  - Slow start, congestion avoidance, fast recovery



additively increase window size ...

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

**cwnd**: TCP sender congestion window size

time

# TCP Congestion Control: Two modes

❖ TCP congestion control is governed by two parameters:

  ▪ **Congestion Window (cwnd)**                    exponential   *2
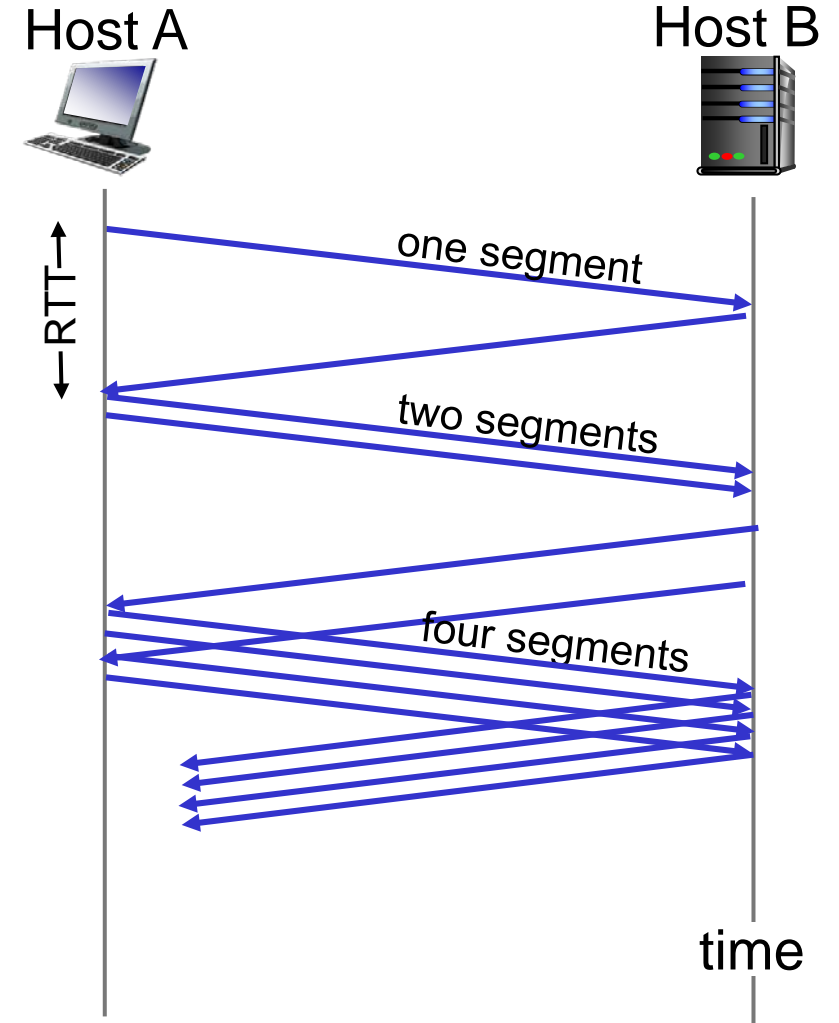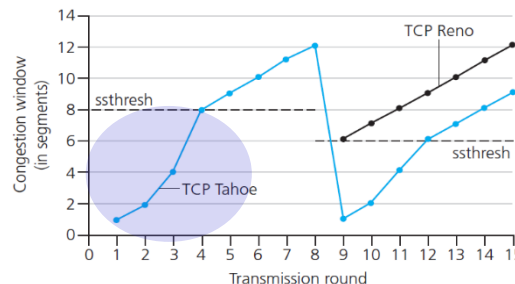                                                    linear
  ▪ **Slow-start threshold Value (ssthresh)**

                                      (slow)              2

❖ Congestion control works in <u>two modes</u>:

  ▪ **Slow Start** (cwnd < ssthresh)
  ▪ **Congestion Avoidance** (cwnd ≥ ssthresh)
  ▪ **(Fast Recovery)**

# TCP "Slow" Start

❖ **when connection begins, increase rate exponentially until first loss event:**
- initially `cwnd` = 1 MSS
- double `cwnd` every RTT
- done by incrementing `cwnd` for every ACK received

❖ *summary:* initial rate is slow but ramps up exponentially fast

# From Slow Start to Congestion Avoidance

TCP slows down the increase of `cwnd` if `cwnd` reaches the slow-start threshold value `ssthresh`

- If cwnd ≥ ssthresh then each time an ACK is received, increment cwnd as follows:
  - cwnd = cwnd + (1 MSS)/cwnd
  (linear increase per RTT)

1

4  -> 1/4 1/4 1/4 1/4
1/4
1

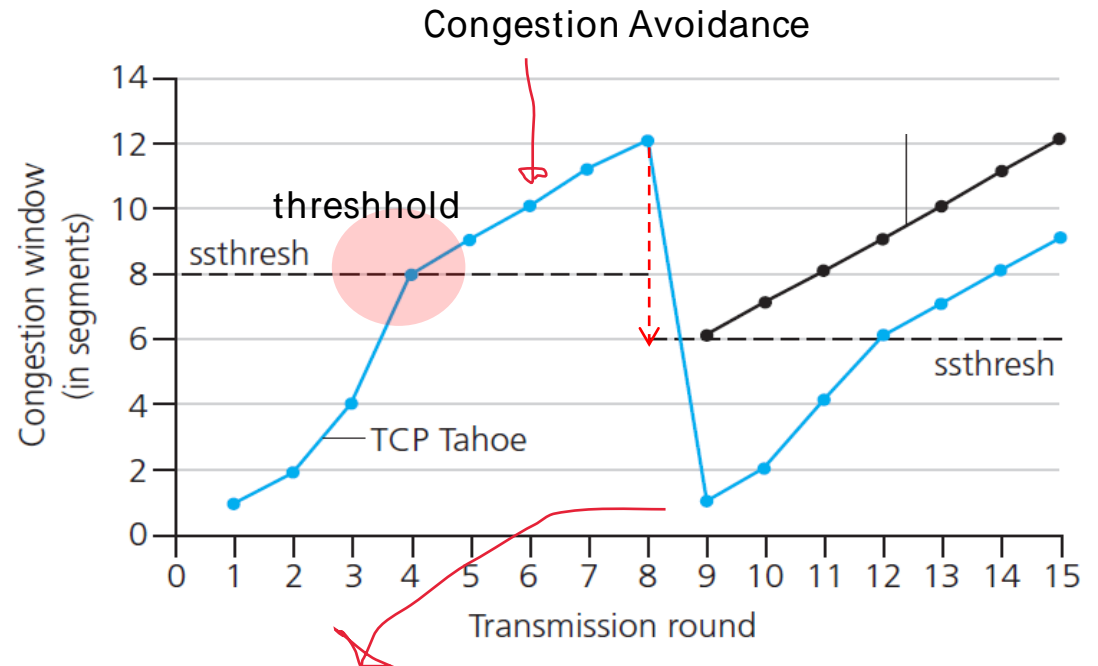## Setting of ssthresh:

- Initially, variable `ssthresh` is set to Advertised window size
- on loss event, `ssthresh` is set to 1/2 of `cwnd` just before loss event

미래창조과학부 SW중심대학
가천대학교 소프트웨어학과
Gachon University Department of Software

# 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 before timeout. = *ssthresh*

# TCP Tahoe Algorithm

**Initially:**
    cwnd = 1 MSS;
    ssthresh = advertised window size;
**New Ack received:**
    if (cwnd < ssthresh)
        /* Slow Start*/
        cwnd = cwnd + 1 MSS;    => 2
    else
        /* Congestion Avoidance */
        cwnd = cwnd + (1 MSS)/cwnd;    => 1
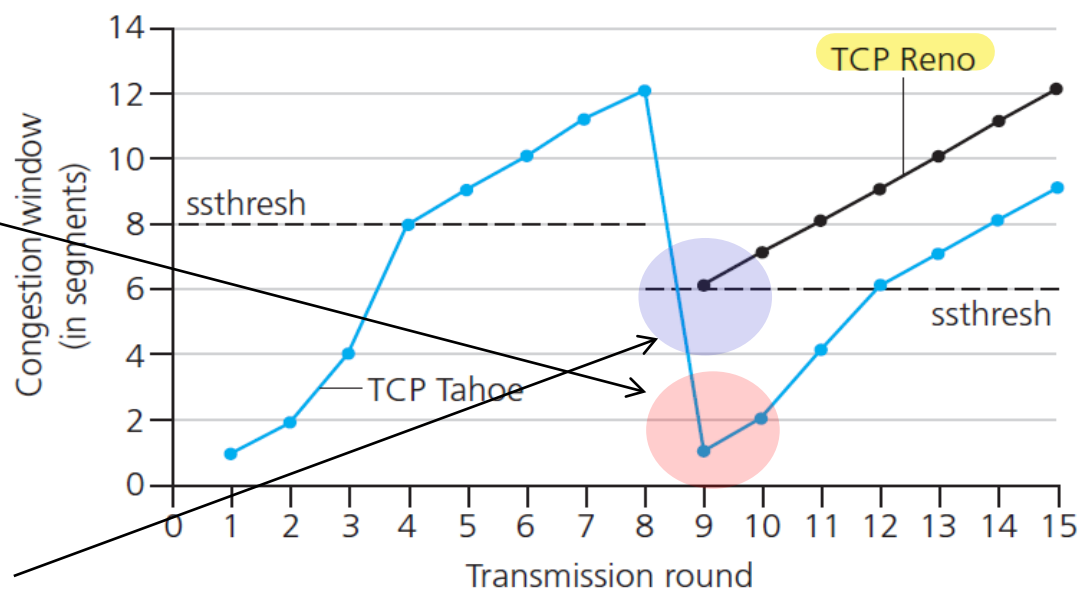**Timeout:**
    /* Multiplicative decrease */
    ssthresh = cwnd/2;
    cwnd = 1 MSS;

# TCP: detecting, reacting to loss

❖ loss indicated by **timeout**:

- ▪ `cwnd` set to 1 MSS;
- ▪ window then grows exponentially (as in slow start) to threshold, then grows linearly

❖ **Fast Recovery (TCP Reno):** loss indicated by 3 duplicate ACKs

- ▪ dup ACKs indicate network capable of delivering some segments (better than timeout)
- ▪ `cwnd` is cut in half window then grows linearly
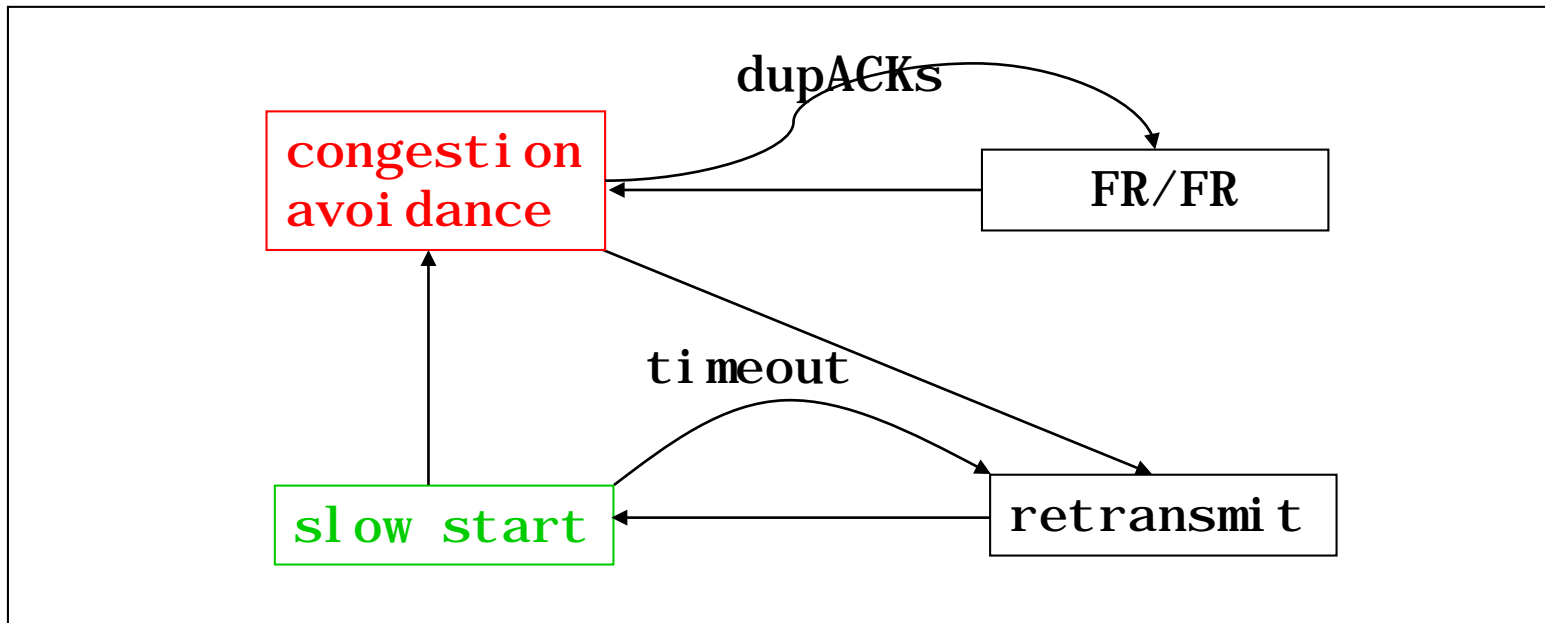


timeout          loss          threshold

# TCP Reno Summary

❖ Basic ideas
  - Fast recovery avoids slow start
  - dupACKs: fast retransmit + fast recovery
  - Timeout: fast retransmit + slow start
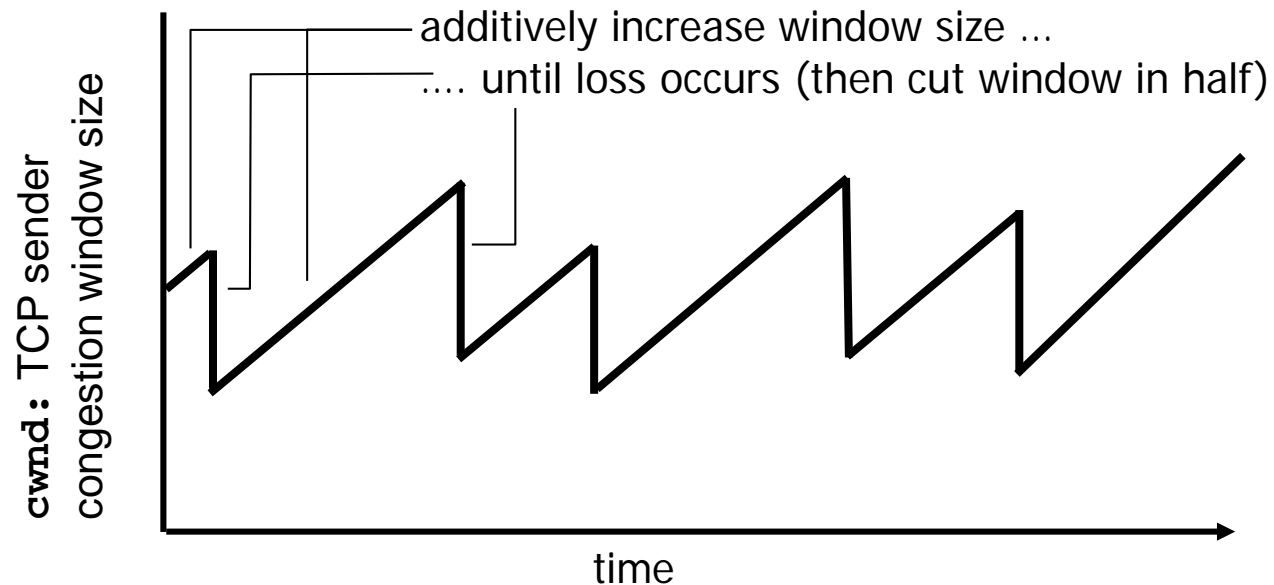
# Additive Increase Multiplicative Decrease (AIMD)

❖ *approach:* sender increases transmission rate (window size), probing for usable bandwidth, until loss occurs

■ *additive increase:* increase `cwnd` by 1 MSS every RTT until loss detected

■ *multiplicative decrease*: cut `cwnd` in half after loss (3-dup acks)

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

미래창조과학부 SW중심대학
가천대학교 소프트웨어학과
Gachon University Department of Software

# Chapter 3: summary

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

미래창조과학부 SW중심대학
가천대학교 소프트웨어학과
Gachon University Department of Software