

# Transport Layer: TCP 2

Oct 15, 2024

Min Suk Kang  
Associate Professor  
School of Computing/Graduate School of Information Security



# Chapter 3: roadmap

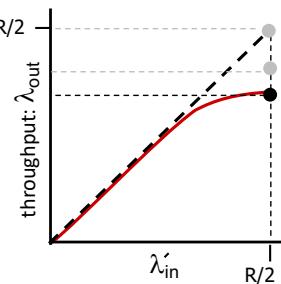
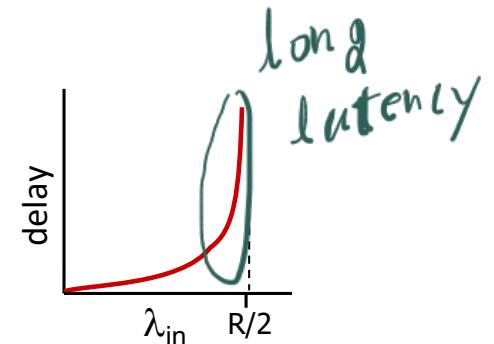
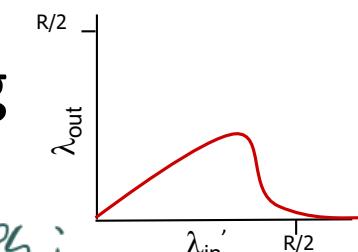
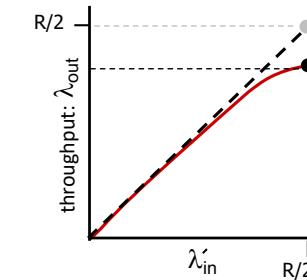
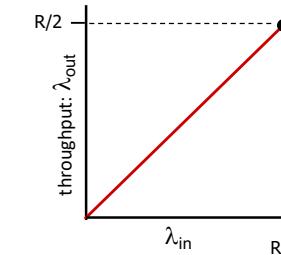
- Transport-layer services
- Multiplexing and demultiplexing
- Connectionless transport: UDP
- Principles of reliable data transfer
- Connection-oriented transport: TCP
- **Principles of congestion control**
- TCP congestion control
- Evolution of transport-layer functionality



# Causes/costs of congestion: insights

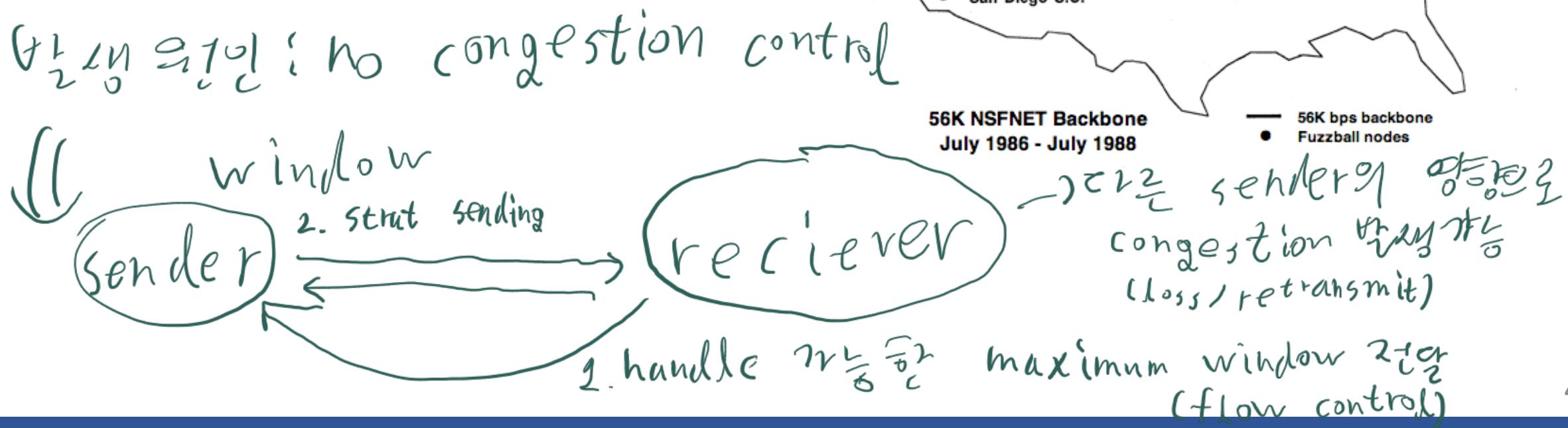
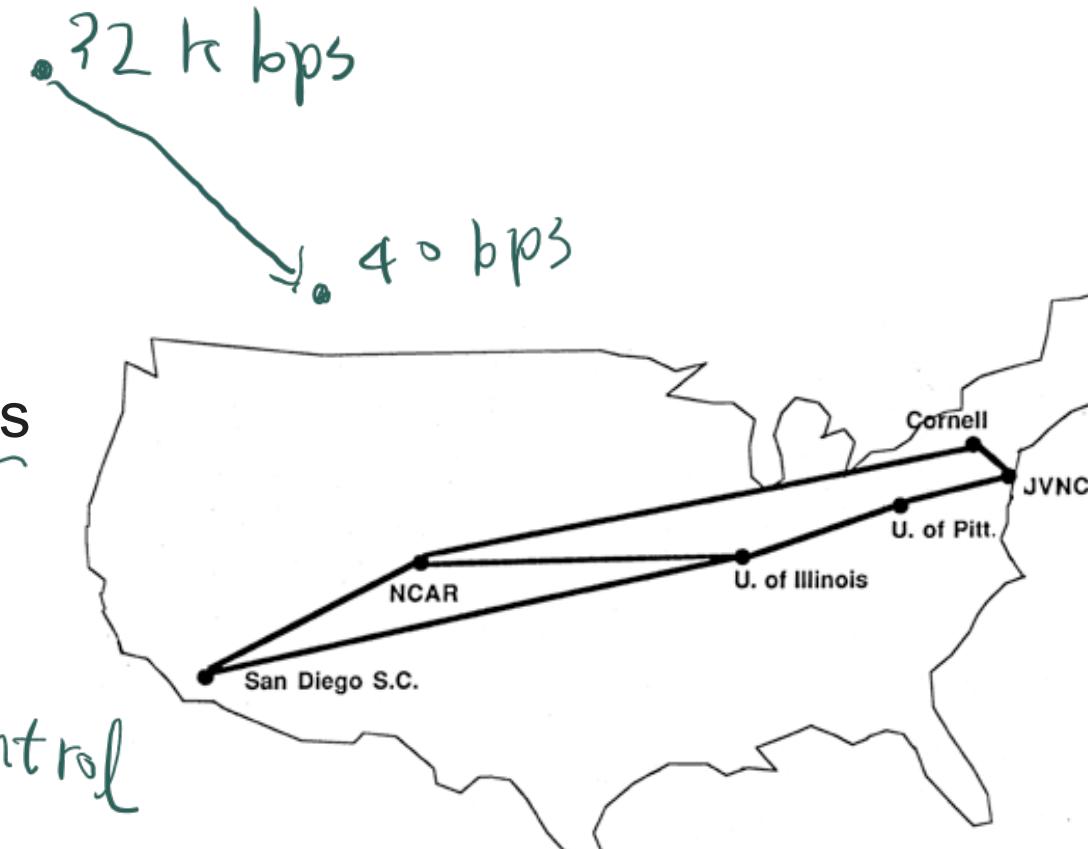
= principle of congestion

- throughput can never exceed capacity
- delay increases as capacity approached
- loss/retransmission decreases effective throughput
  - ↳ by time out
- un-needed duplicates further decreases effective throughput
- upstream transmission capacity / buffering wasted for packets lost downstream
  - ↳ wasted all sources in upstream



# “Congestion Collapse”

- In October 1986, NSFNET backbone capacity dropped from 32 kbps to 40 bps
- TCP congestion control invented between 1987-1988 by Van Jacobson and Sally Floyd

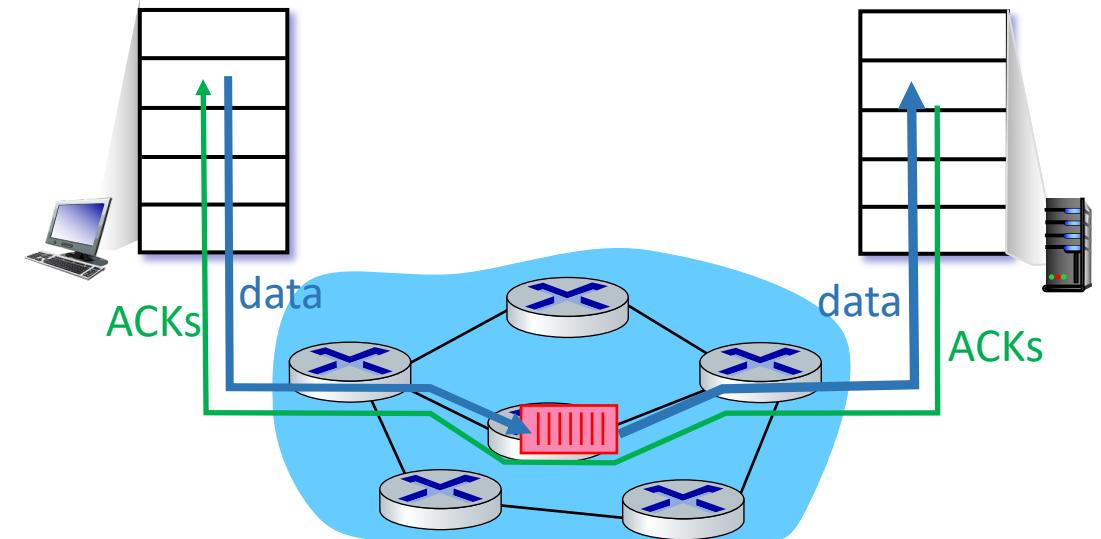


# Approaches towards congestion control

## End-end congestion control:

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

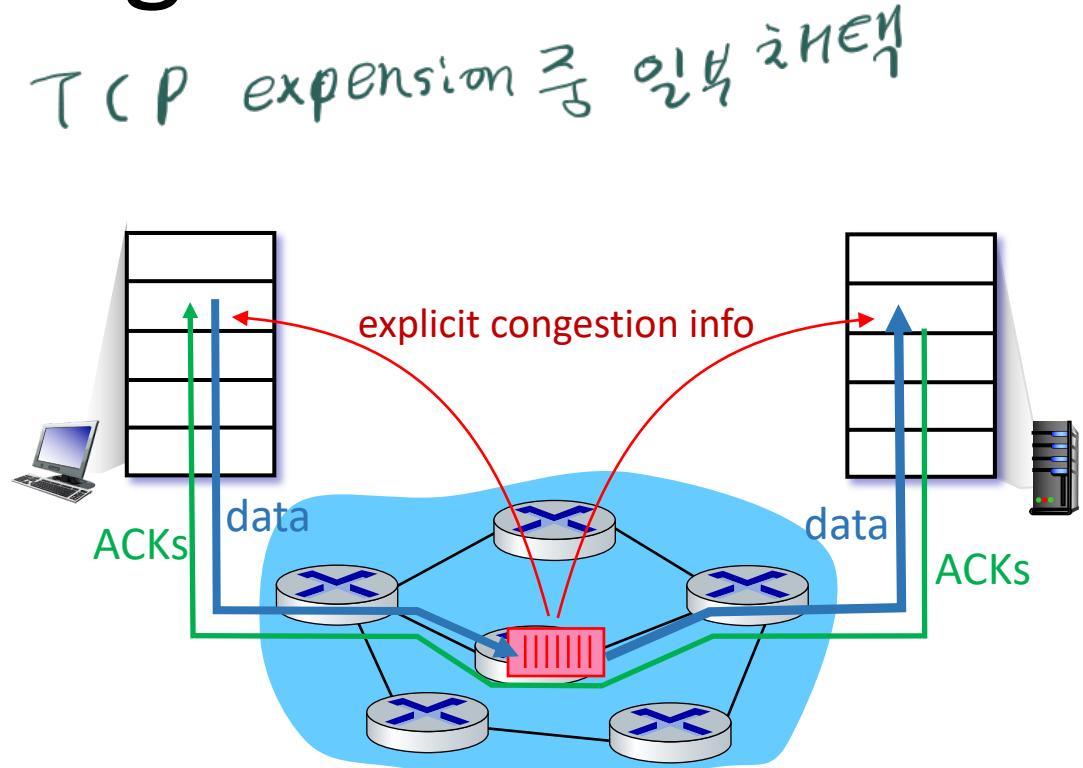
loss, delay  $\frac{2}{2}$   $\frac{5}{6} \leq H$  congestion  $\frac{2}{2} \frac{2}{7}$



# Approaches towards congestion control

## Network-assisted congestion control:

- routers provide direct feedback to sending/receiving hosts with flows passing through congested router
- may indicate congestion level or explicitly set sending rate
- TCP ECN, ATM, DECbit protocols



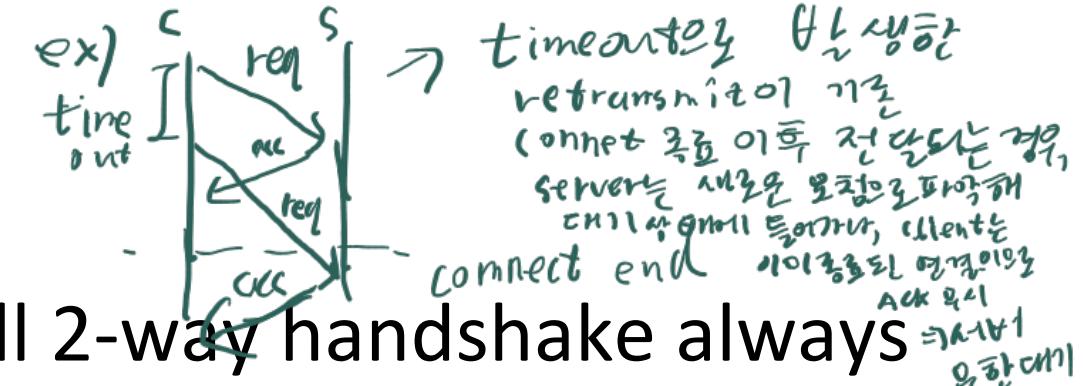
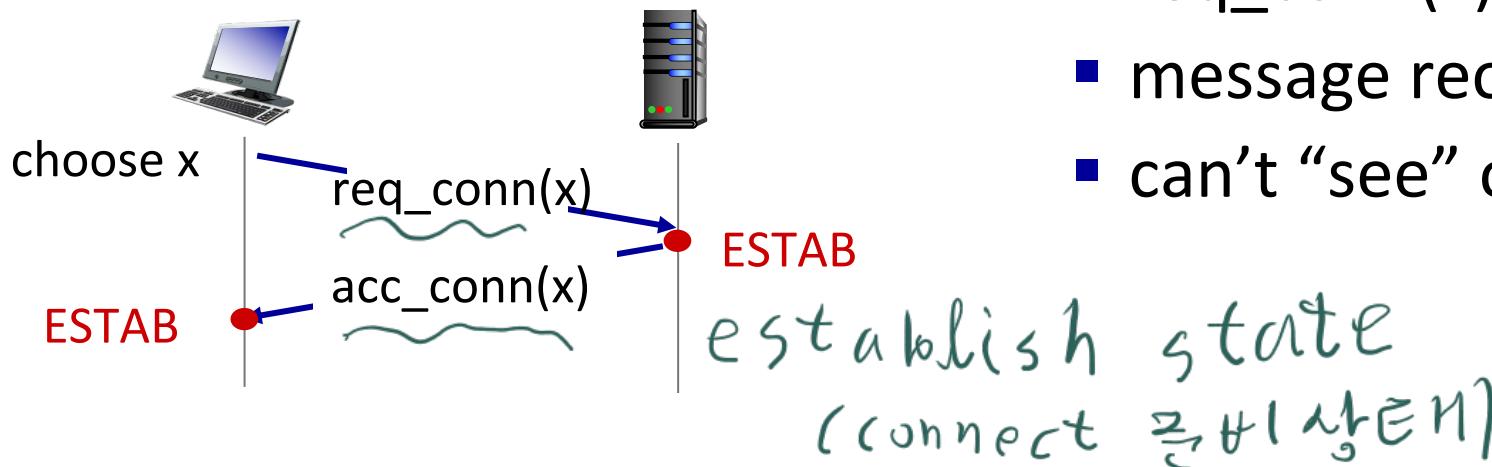
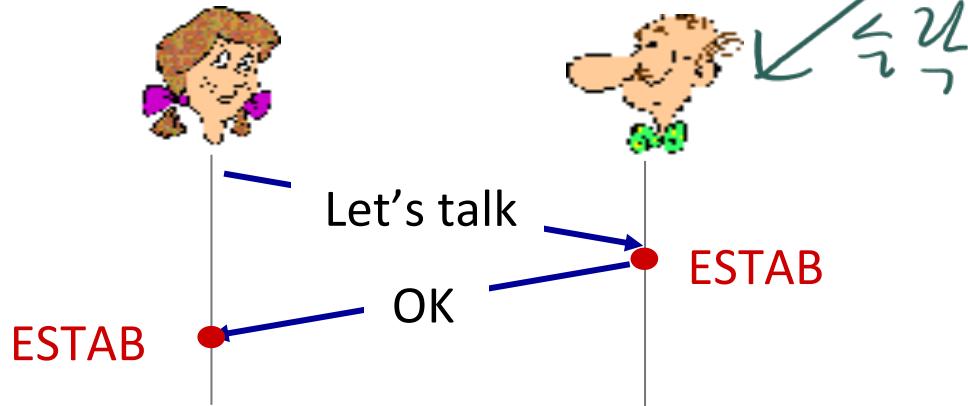
# Chapter 3: roadmap

- Transport-layer services
- Multiplexing and demultiplexing
- Connectionless transport: UDP
- Principles of reliable data transfer
- **Connection-oriented transport: TCP**
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- Principles of congestion control
- TCP congestion control



# Agreeing to establish a connection

2-way handshake:



Q: will 2-way handshake always work in network?

- variable delays → WT delay
- retransmitted messages (e.g. req\_conn(x)) due to message loss → ok drop
- message reordering → segment 순서 혼동
- can't "see" other side → 직접 보

State 초기화  
복구

# TCP 3-way handshake

## Client state

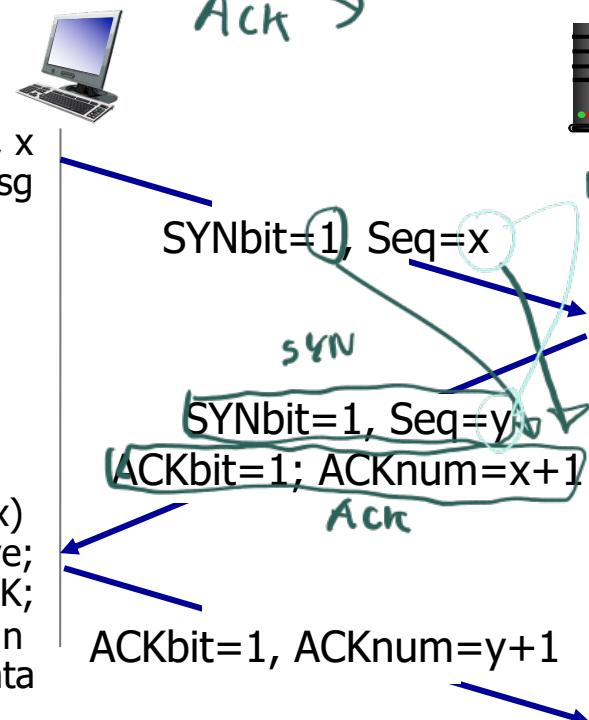
```
clientSocket = socket(AF_INET, SOCK_STREAM)
```

LISTEN

```
clientSocket.connect((serverName, serverPort))
```

SYNSENT

choose init seq num, x  
send TCP SYN msg



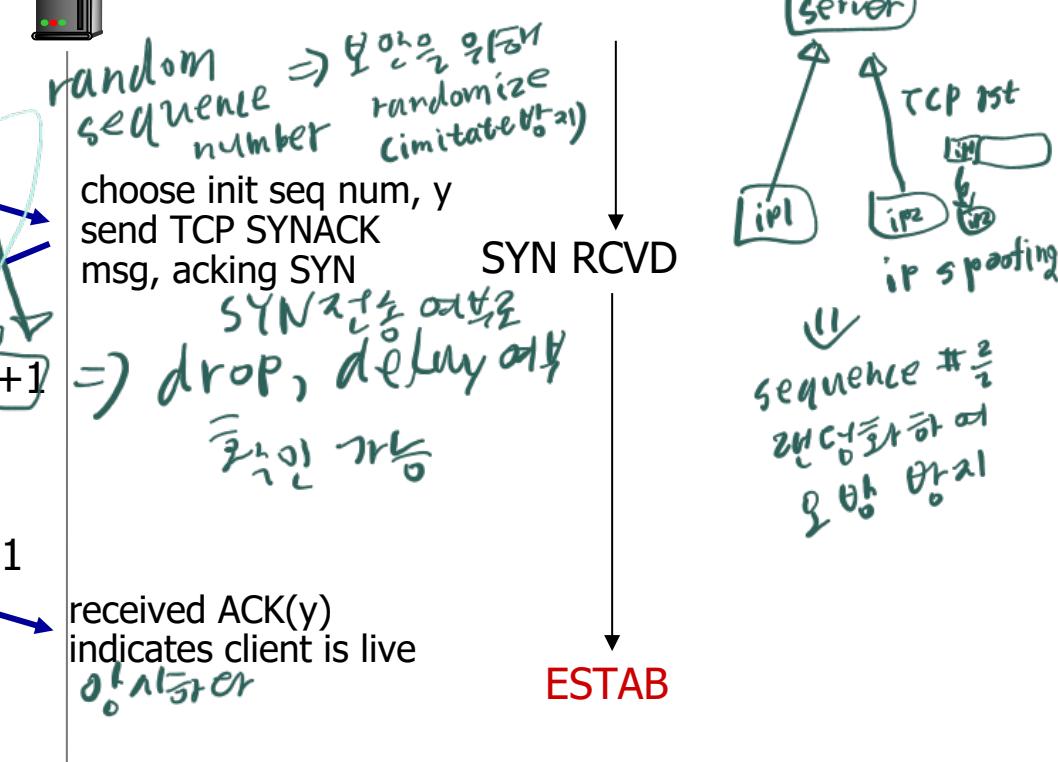
ESTAB

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

## Server state

```
serverSocket = socket(AF_INET, SOCK_STREAM)
serverSocket.bind(('', serverPort))
serverSocket.listen(1)
connectionSocket, addr = serverSocket.accept()
```

LISTEN



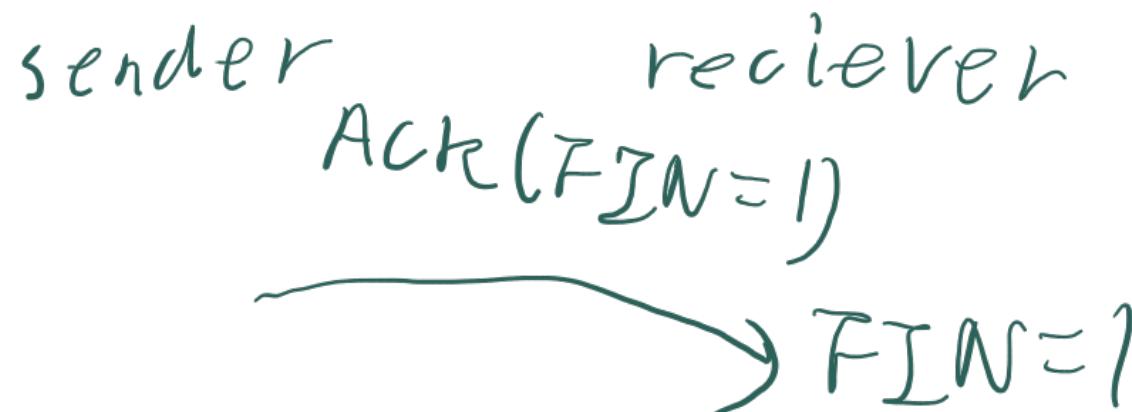
# A human 3-way handshake protocol



# Closing a TCP connection

- client, server each close their side of connection
  - send TCP segment with FIN bit = 1
- respond to received FIN with ACK
  - on receiving FIN, ACK can be combined with own FIN
- simultaneous FIN exchanges can be handled

도시락



# Chapter 3: roadmap

- Transport-layer services
- Multiplexing and demultiplexing
- Connectionless transport: UDP
- Principles of reliable data transfer
- Connection-oriented transport: TCP
- Principles of congestion control
- **TCP congestion control**
- Evolution of transport-layer functionality



# Thought experiment: subway chair



7가지의 조건  
사람이 앉을  
때는 7개를  
전부 차지해도  
OK  
문제는 다른  
사람과 resource<sup>2</sup>을  
공유해야 한다.  
①  
즉 어떤 상황을 갖는  
업으로, TCP는 자원  
낭비도 하지 않자고  
⇒ 효율 높일 방법?  
일단 예고, 복습하면 13

# TCP congestion control: AIMD

- **approach:** senders can increase sending rate until packet loss (congestion) occurs, then decrease sending rate on loss event

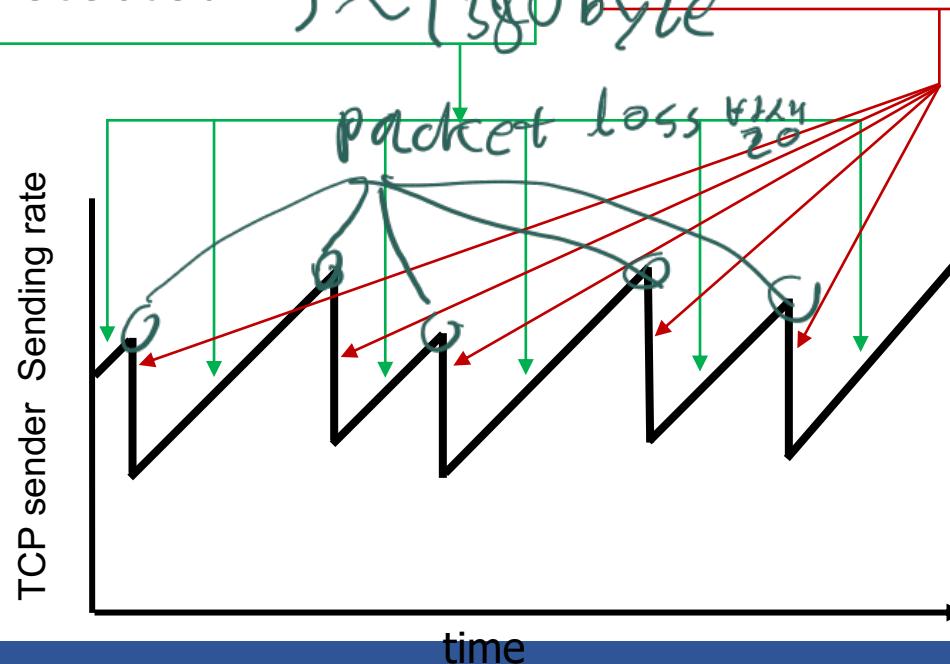
Congestion avoidance

## - Additive Increase

increase sending rate by 1  
maximum segment size every  
RTT until loss detected

## Multiplicative Decrease

cut sending rate in half at each loss event



**AIMD** sawtooth behavior: probing for bandwidth

# TCP AIMD: more

**Multiplicative decrease** detail: sending rate is

- Cut in half on loss detected by triple duplicate ACK (TCP Reno)
- Cut to 1 MSS (maximum segment size) when loss detected by timeout (TCP Tahoe)

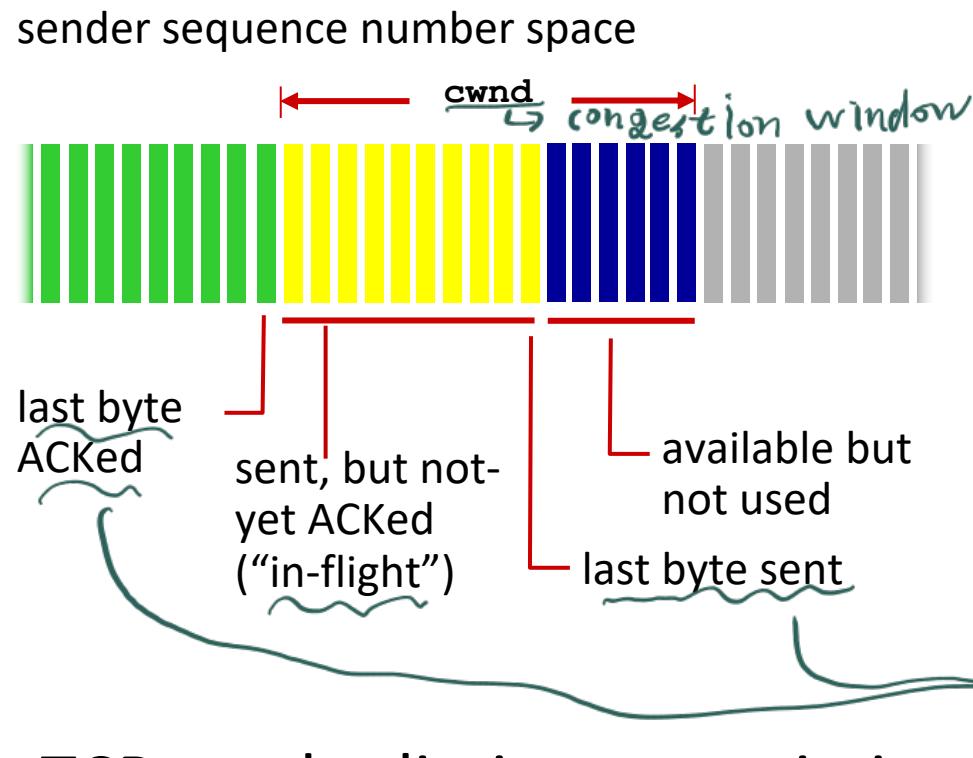
3 A를 중복 시킬 때  $\times \frac{1}{2}$   
(TCP fast retransmission)  
→ major problem 발생 시 1초로 커짐  
 $(= \text{timeout})$   
1 MSS 단위 = packet 1M

Why AIMD?

- AIMD – a distributed, asynchronous algorithm – has been shown to: 부상도 비동기
  - optimize congested flow rates network wide!
  - have desirable stability properties

안정성

# TCP congestion control: details



- TCP sender limits transmission:  $\text{LastByteSent} - \text{LastByteAcked} \leq \text{cwnd}$
- cwnd is dynamically adjusted in response to observed network congestion (implementing TCP congestion control)

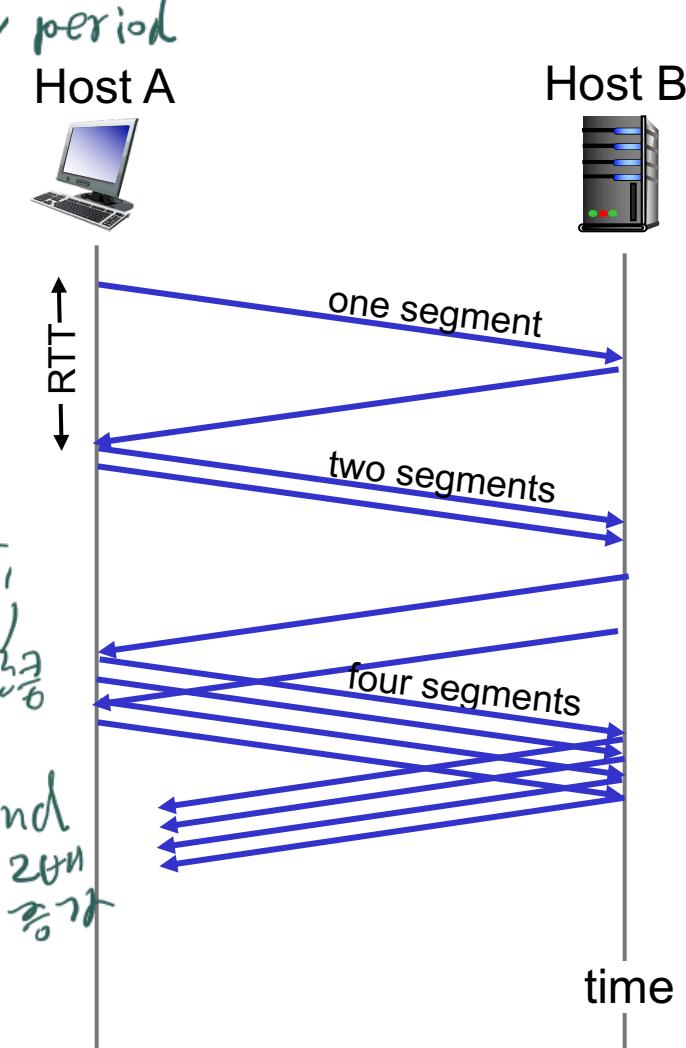
↳ cwnd를 congestion 유통에 따라 동적 조정 {  
1. slow start  
2. congestion avoidance  
3. fast recovery}

# TCP "slow" start

↳ compare to 'no-congestion control' period

- 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

혼잡 통제가 업신연 시설에 의해 slow, 실제로는 fast  
(초기부터 최대값으로 증가)



# TCP: from slow start to congestion avoidance

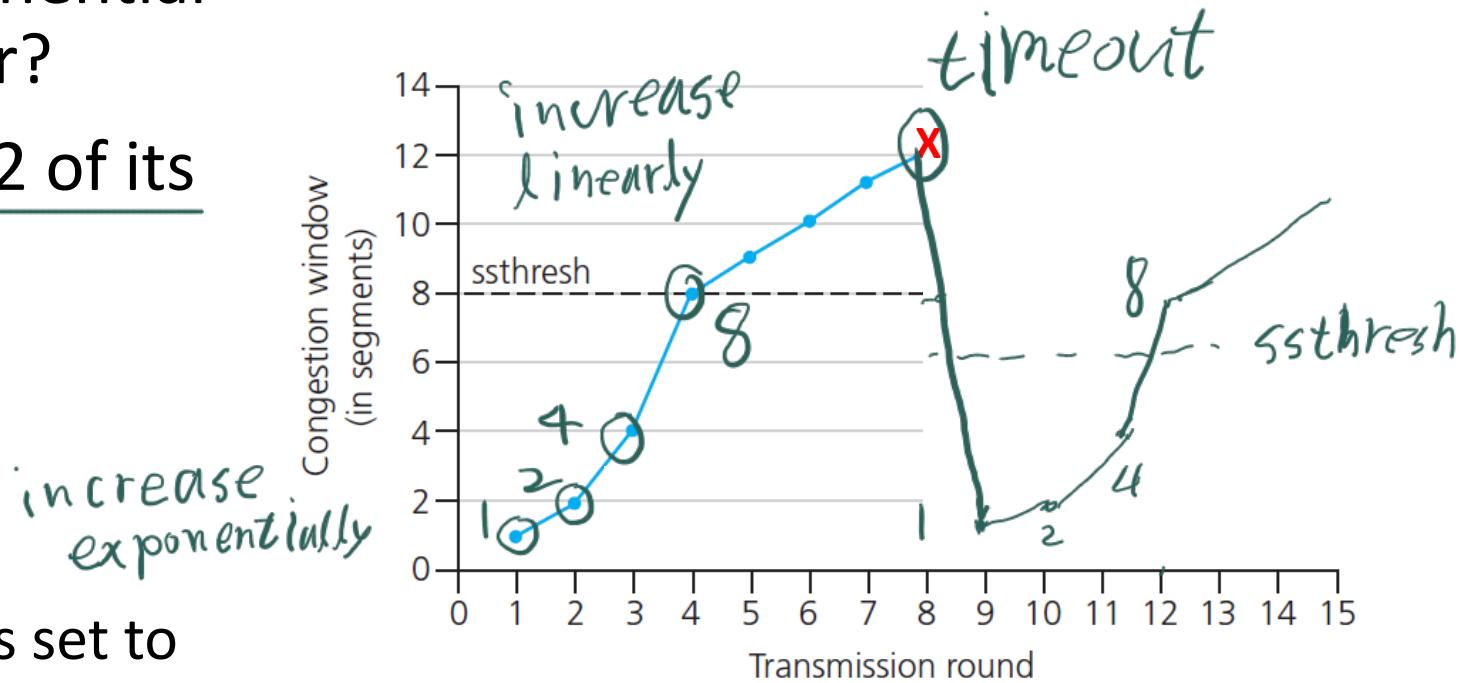
**Q:** when should the exponential increase switch to linear?

**A:** when cwnd gets to  $1/2$  of its value before timeout.

## Implementation:

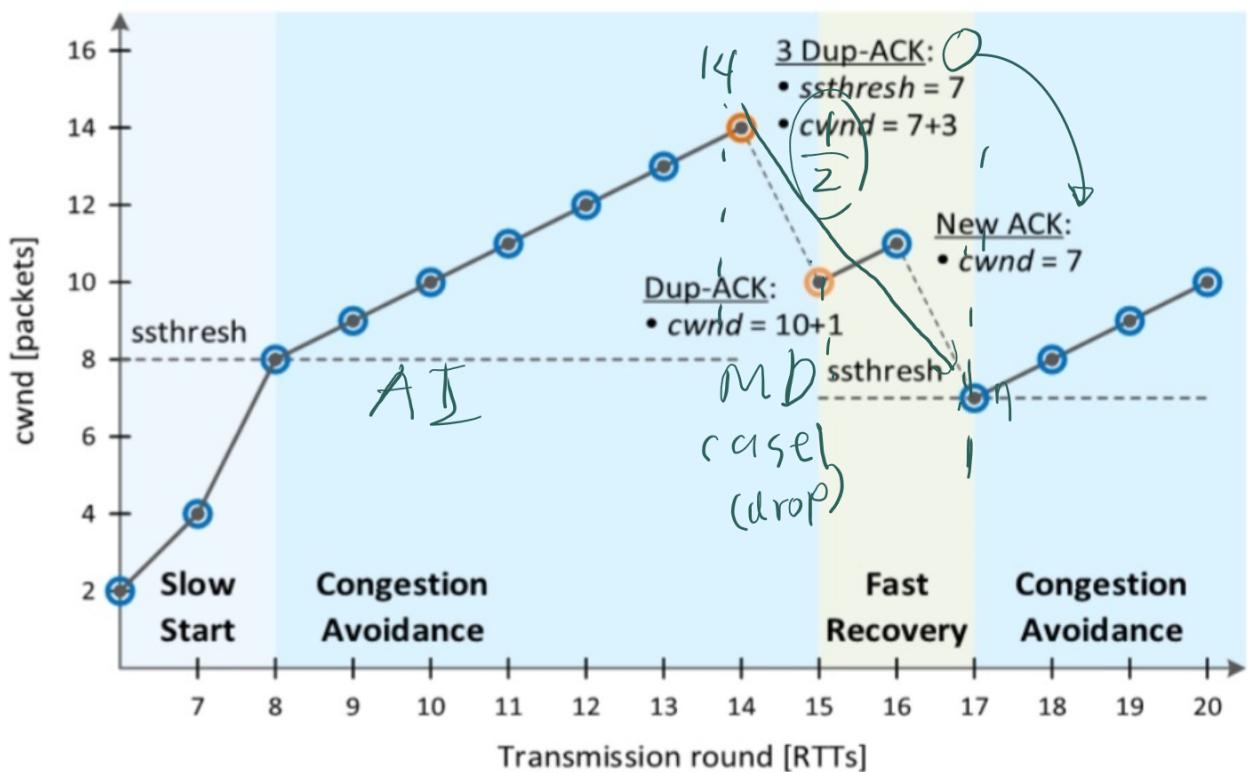
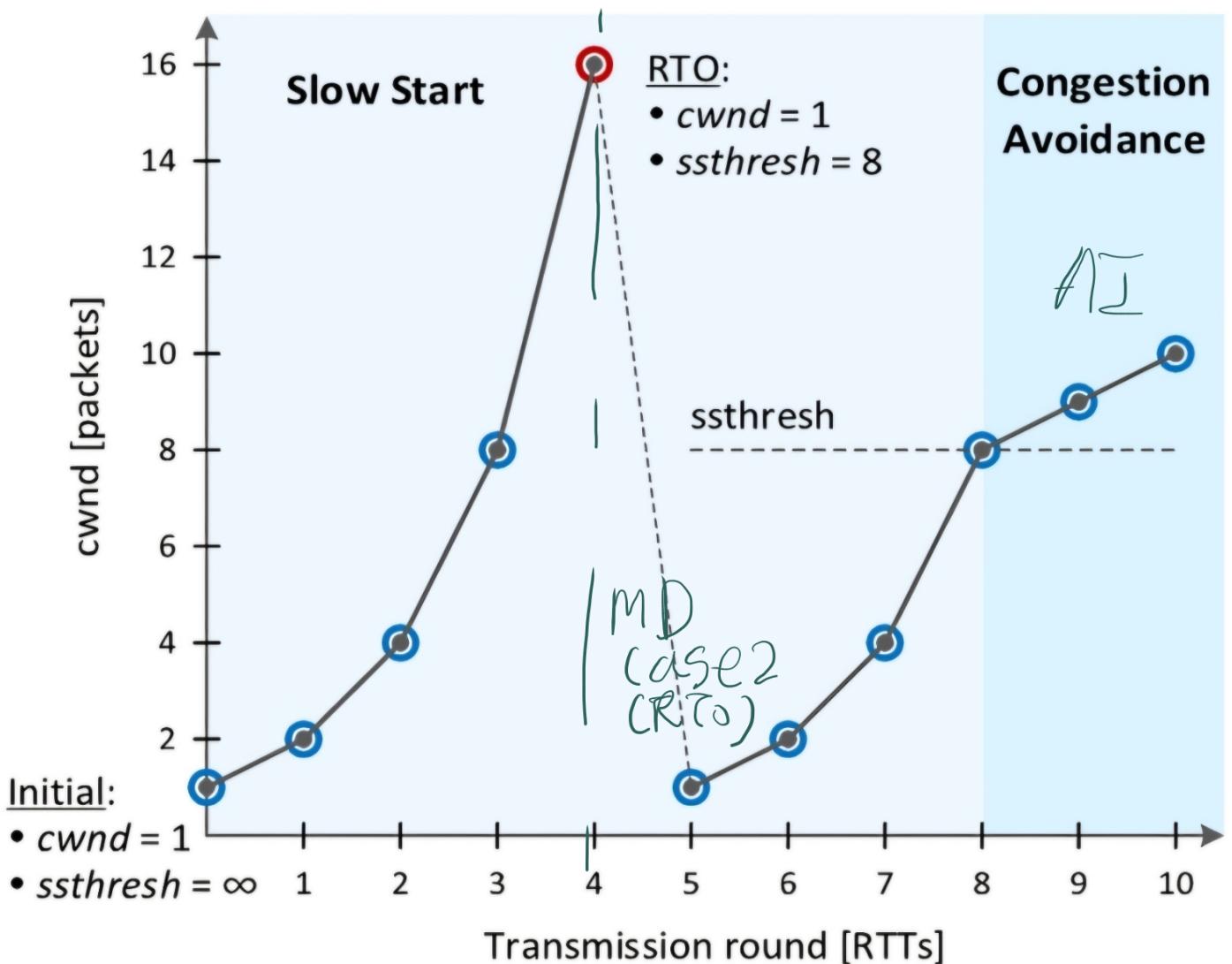
- variable ssthresh
- on loss event, ssthresh is set to  $1/2$  of cwnd just before loss event

packet drop  
발생으로 인한  
RTO (retransmission timeout)

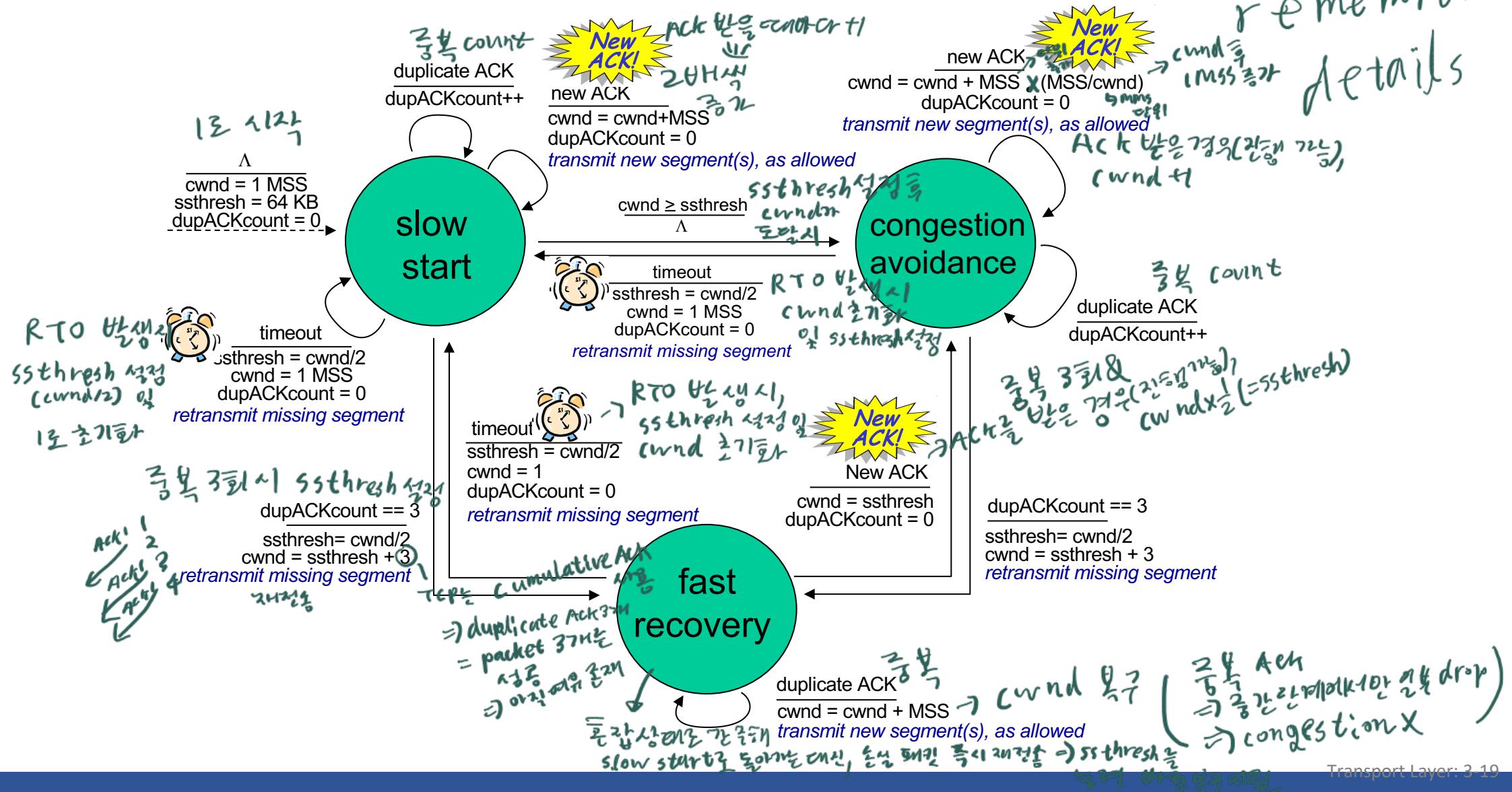


직전값의  $\frac{1}{2}$ 로 ssthresh를 12로 초기화  
cwnd는 12로 초기화

\* Check out the online interactive exercises for more examples: [http://gaia.cs.umass.edu/kurose\\_ross/interactive/](http://gaia.cs.umass.edu/kurose_ross/interactive/)



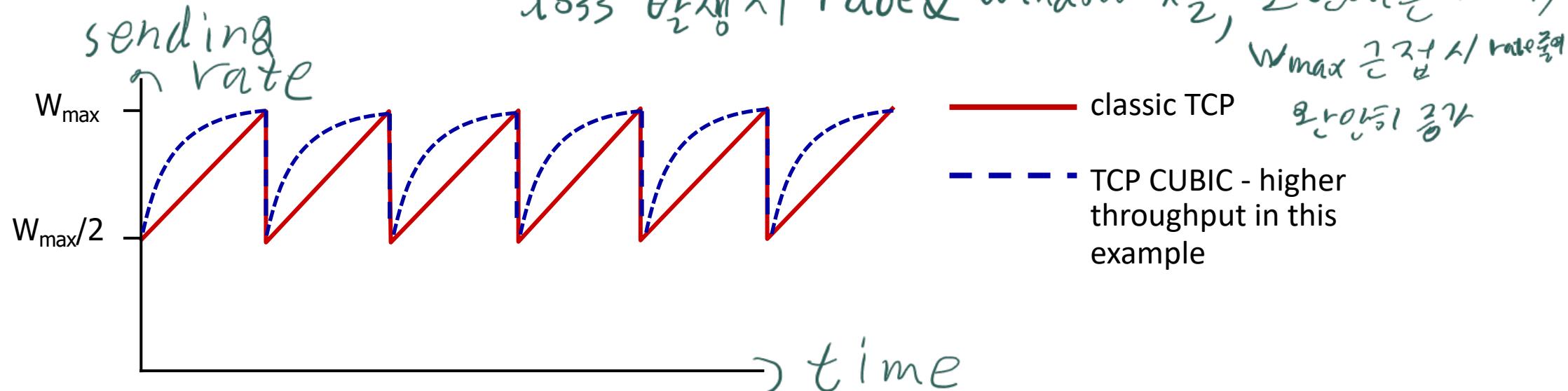
# Summary: TCP congestion control



# TCP CUBIC

- Is there a better way than AIMD to “probe” for usable bandwidth?
- Insight/intuition:

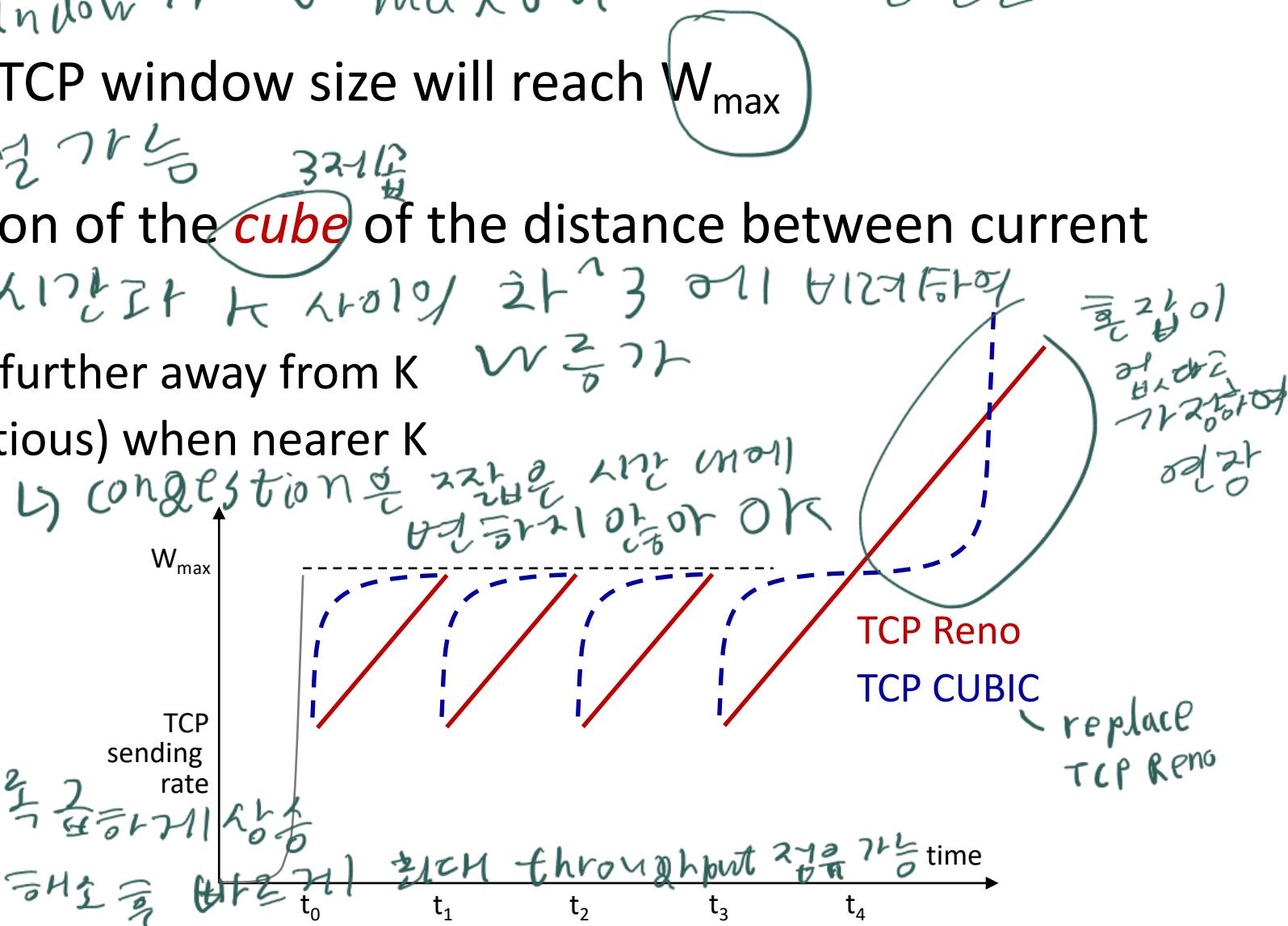
- $W_{\max}$ : sending rate at which congestion loss was detected
- congestion state of bottleneck link probably (?) hasn't changed much
- after cutting rate/window in half on loss, initially ramp to  $W_{\max}$  faster, but then approach  $W_{\max}$  more slowly



# TCP CUBIC

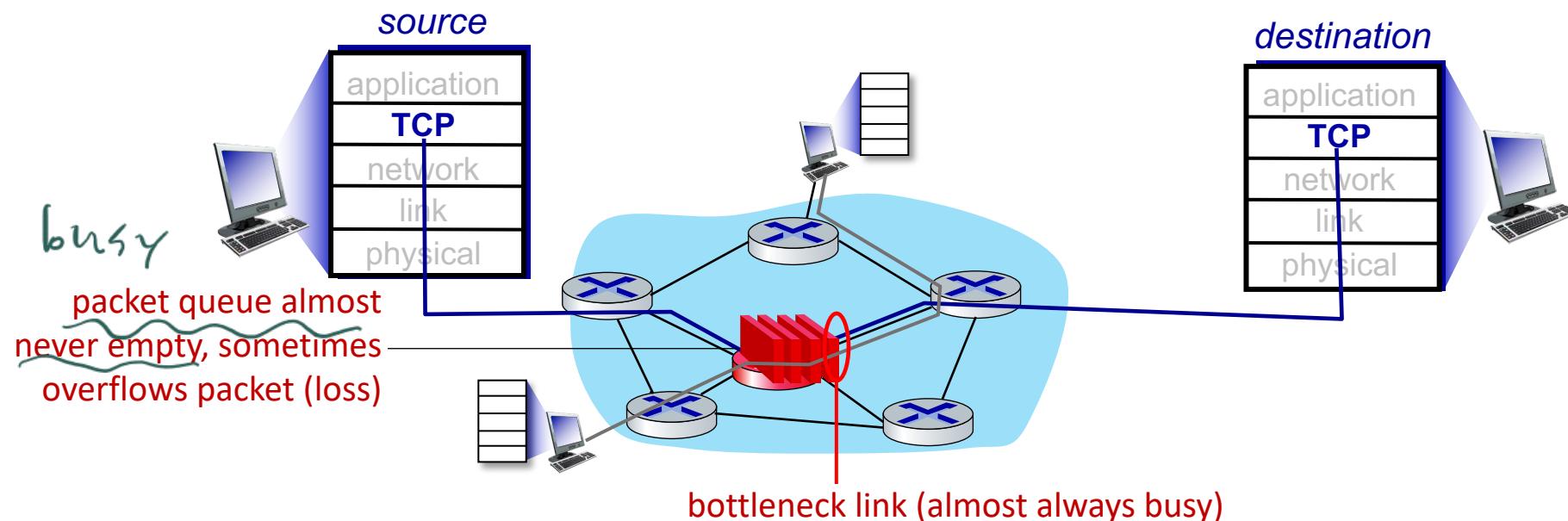
window가  $W_{max}$ 에 가까울 때 순간

- K: point in time when TCP window size will reach  $W_{max}$ 
  - K itself is tunable 조절 가능 32bit
- increase W as a function of the cube of the distance between current time and K 현재 시간과 K 사이의 차  $^3$ 로 A인 경우  
$$w_{(rate)} = (k-t)^{*3}$$
  - larger increases when further away from K  $w \frac{2}{3}$  가
  - smaller increases (cautious) when nearer K
- TCP CUBIC default in Linux, most popular TCP for popular Web servers



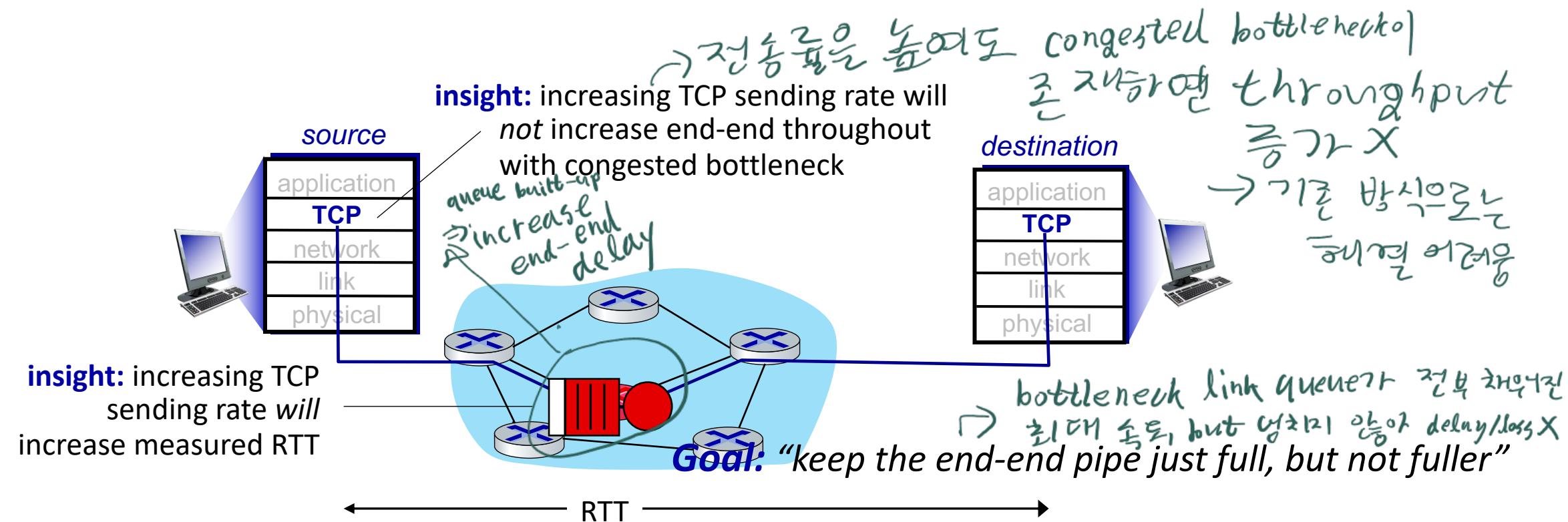
# TCP and the congested “bottleneck link”

- TCP (classic, CUBIC) increase TCP's sending rate until packet loss occurs at some router's output: the *bottleneck link*



# TCP and the congested “bottleneck link”

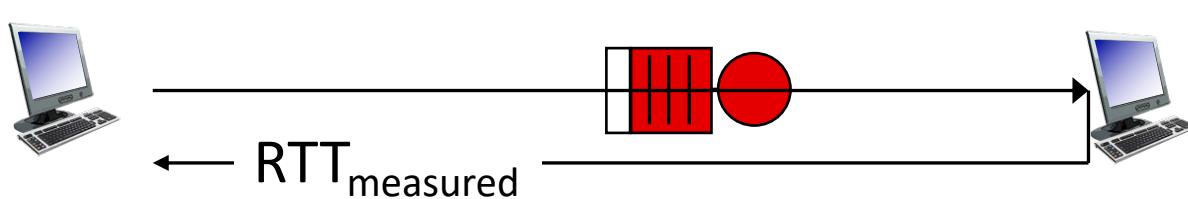
- TCP (classic, CUBIC) increase TCP's sending rate until packet loss occurs at some router's output: the *bottleneck link*
- understanding congestion: useful to focus on congested bottleneck link



# Delay-based TCP congestion control

↪ bottleneck link에 집중하는 방법

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



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

## Delay-based approach:

- $\text{RTT}_{\min}$  - minimum observed RTT (uncongested path) 혼잡 X
- uncongested throughput with congestion window  $cwnd$  is  $cwnd/\text{RTT}_{\min}$ 
  - if measured throughput “very close” to uncongested throughput → 측정값과 최대치 유포  
increase  $cwnd$  linearly /\* since path not congested \*/ ⇒ 혼잡 X 이므로  $cwnd↑$
  - else if measured throughput “far below” uncongested throughput → 측정값과 최대치 차이↑  
decrease  $cwnd$  linearly /\* since path is congested \*/ ⇒ 혼잡 상태이므로  $cwnd↓$

# Delay-based TCP congestion control

- congestion control without inducing/forcing loss 유도  
loss 풀 땅 차지율  
로스率 ⇒ loss 풀 땅  
업로드 흐름  
통제
- maximizing throughout ("keeping the just pipe full...") while keeping delay low ("...but not fuller")
- a number of deployed TCPs take a delay-based approach
  - BBR deployed on Google's (internal) backbone network

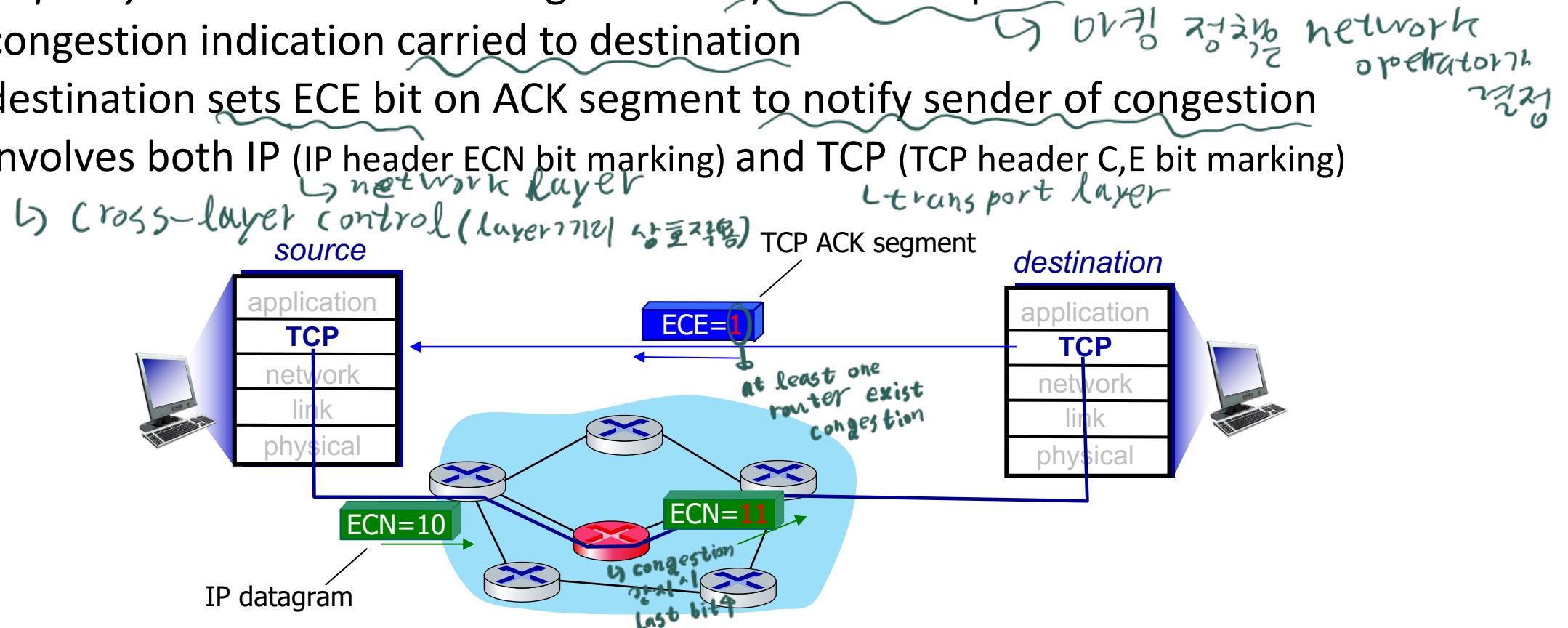


→ BBR의 backbone에서만  
사용되는 이유  
⇒ public networks  
작은 RTT와 사용자 앱  
경로 고정, 고정된으로 RTT min 확장  
traffic 정적 → 입력과 RTT 및 bottleneck  
hand width이 가정에 따라  
prediction 오류

# Explicit congestion notification (ECN)

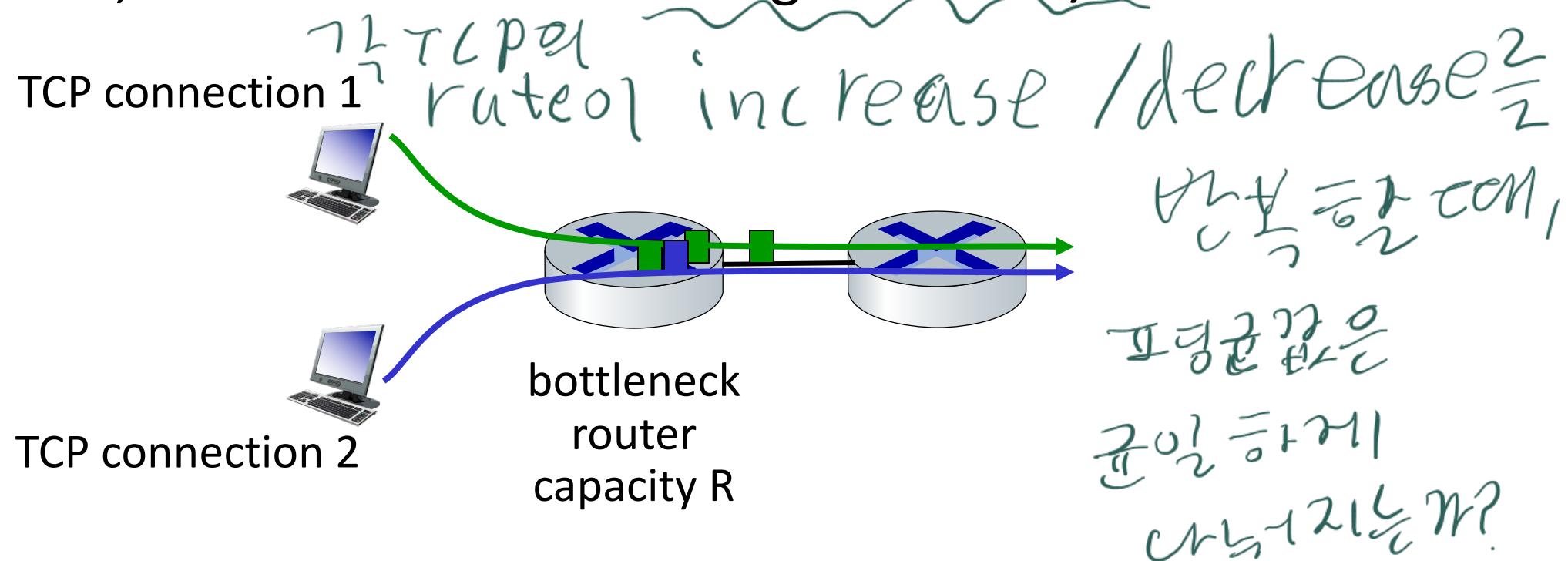
TCP deployments often implement network-assisted congestion control:

- two bits in IP header (ToS field) marked by network router to indicate congestion
  - *policy to determine marking chosen by network operator*
- congestion indication carried to destination
- destination sets ECE bit on ACK segment to notify sender of congestion
- involves both IP (IP header ECN bit marking) and TCP (TCP header C,E bit marking)



# TCP fairness

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

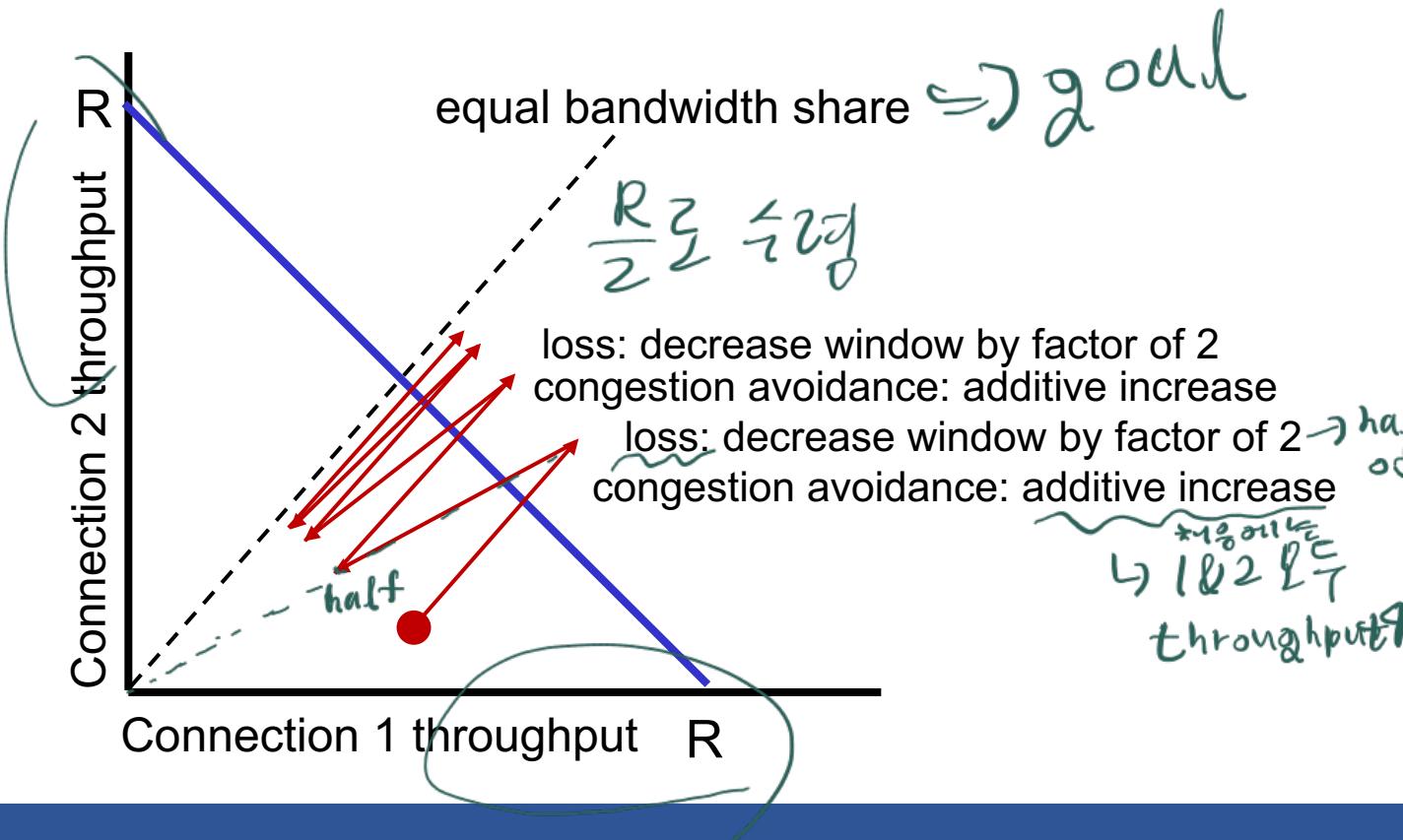


# Q: is TCP Fair?

Example: two competing TCP sessions:

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

→ fair without negotiate



Is TCP fair?

A: Yes, under idealized assumptions:

- same RTT
- fixed number of sessions  
only in congestion avoidance  $\rightarrow$  no added session

# Fairness: must all network apps be “fair”?

→ fairness is “hot” to bust → monitor the violent (TCP를 2로 나누기) traffic  
Fairness and UDP

- multimedia apps often do not use TCP

- do not want rate throttled by congestion control → congestion control X

- instead use UDP:

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

- there is no “Internet police” policing use of congestion control

UDP는 혼잡 통제가

없으나 2가지 사용

## 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

TCP는 connection 단위 fairness

구현할 때 app에서 다른 연결 사용 시 fairness 보장 X

# Transport layer: roadmap

- Transport-layer services
- Multiplexing and demultiplexing
- Connectionless transport: UDP
- Principles of reliable data transfer
- Connection-oriented transport: TCP
- Principles of congestion control
- TCP congestion control
- Evolution of transport-layer functionality



# Evolving transport-layer functionality

- TCP, UDP: principal transport protocols for 40 years
- different “flavors” of TCP developed, for specific scenarios:

Scenario	Challenges
Long, fat pipes (large data transfers)	Many packets “in flight”; loss shuts down pipeline <i>frequent loss</i>
Wireless networks	Loss due to noisy wireless links, mobility; TCP treat this as congestion loss
Long-delay links	Extremely long RTTs
Data center networks	Latency sensitive
Background traffic flows	Low priority, “background” TCP flows

→ TCP چیزی نیست

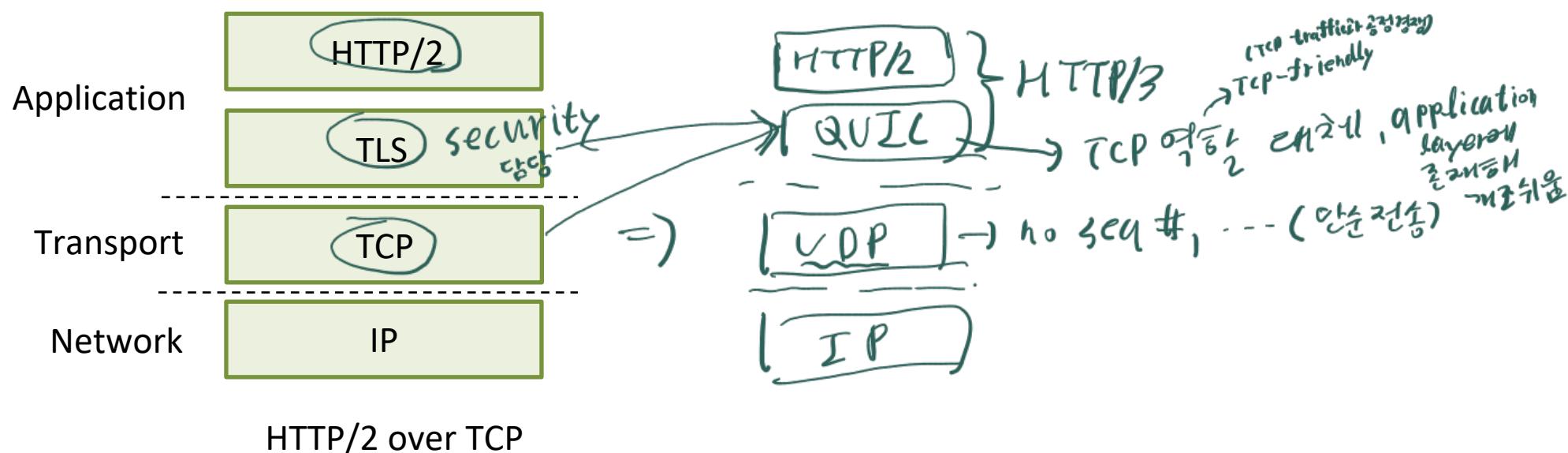
- moving transport-layer functions to application layer, on top of UDP

- HTTP/3: QUIC

چیزی نیست

# QUIC: Quick UDP Internet Connections

- application-layer protocol, on top of UDP
  - increase performance of HTTP
  - deployed on many Google servers, apps (Chrome, mobile YouTube app)



# Mixed feelings about HTTP/3 (QUIC)

- <https://youtu.be/wV9FSyFB8tk>



Horrible, Helpful, http3 Hack - Computerphile  
72K views • 3 months ago

 Computerphile ✓

http3 is here, but it wasn't an easy solution, Richard G Clegg of Queen

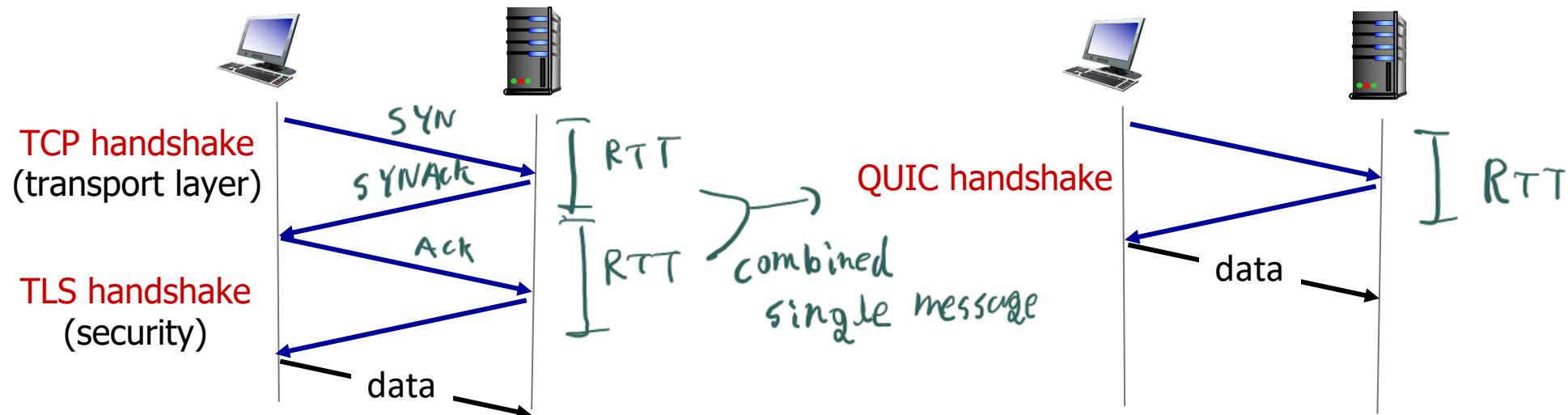
4K

# QUIC: Quick UDP Internet Connections

adopts approaches we've studied in this chapter for connection establishment, error control, congestion control

- error and congestion control: “Readers familiar with TCP’s loss detection and congestion control will find algorithms here that parallel well-known TCP ones.” [from QUIC specification]
- connection establishment: reliability, congestion control, authentication, encryption, state established in one RTT
- multiple application-level “streams” multiplexed over single QUIC connection
  - separate reliable data transfer, security
  - common congestion control

# QUIC: Connection establishment



TCP (reliability, congestion control state) + TLS (authentication, crypto state)

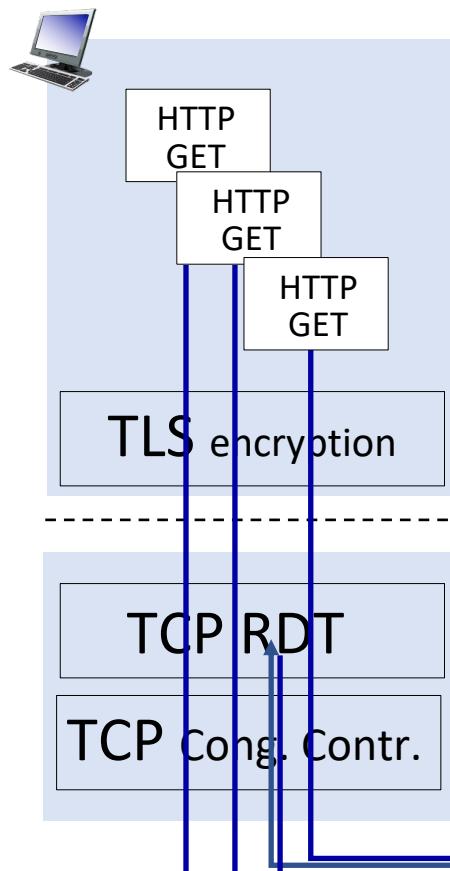
- 2 serial handshakes

QUIC: reliability, congestion control, authentication, crypto state

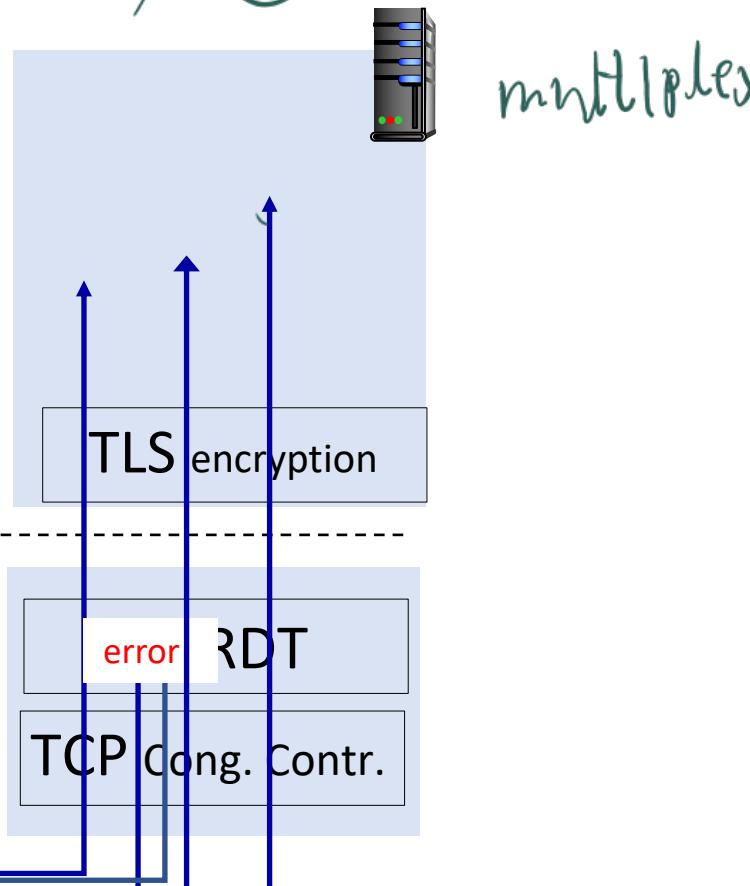
- 1 handshake

# QUIC: streams: parallelism, no HOL blocking

application

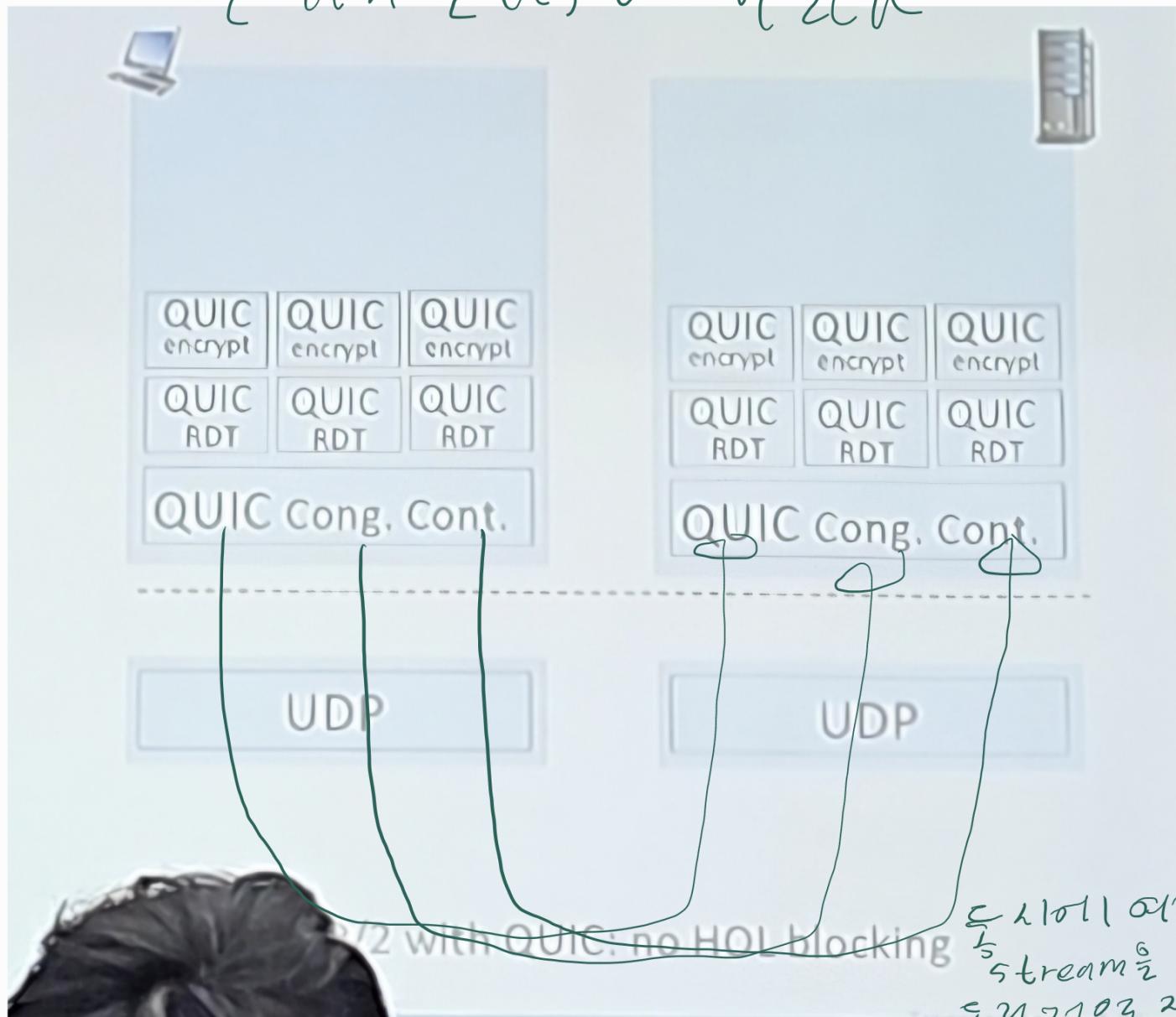


(a) HTTP 1.1



error  
delay

can customized



# Chapter 3: summary

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

## Up next:

- leaving the network “edge” (application, transport layers)
- into the network “core”
- two network-layer chapters:
  - data plane
  - control plane

# Next...

- *4.1 Overview of Network Layer*
  - *4.2 What's Inside a Router?*
- +
- *Review last year's midterm questions*