

Transport Layer: TCP 1

Oct 10, 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**
 - segment structure
 - reliable data transfer
 - flow control
 - connection management
- Principles of congestion control
- TCP congestion control



TCP: overview

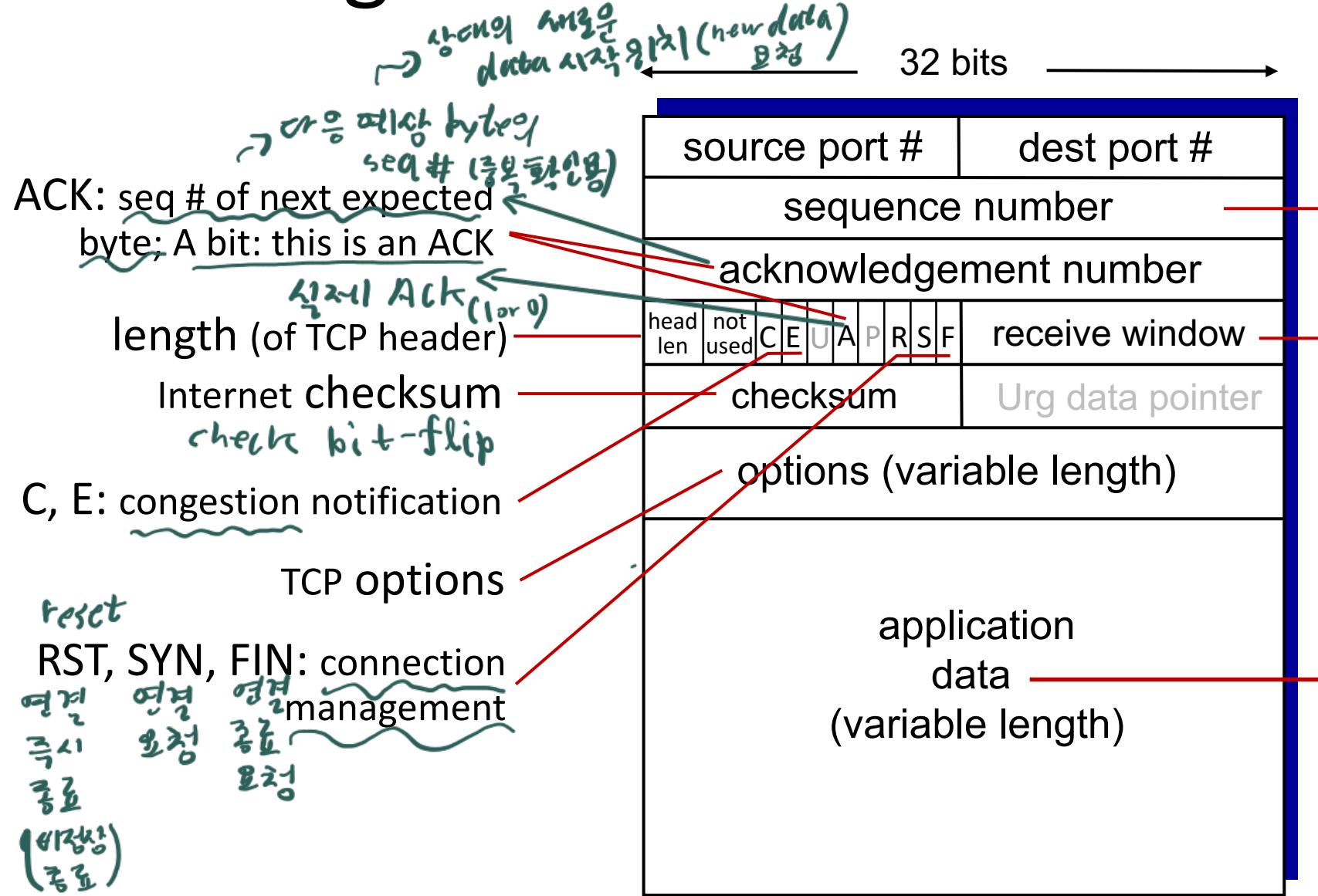
RFCs: 793, 1122, 2018, 5681, 7323

↳ document

Go-back-N

- point-to-point:
 - one sender, one receiver
- reliable, in-order *byte steam*:
 - no "message boundaries"
- full duplex data:
 - bi-directional data flow in same connection
 - MSS: maximum segment size
- cumulative ACKs
- pipelining:
 - TCP congestion and flow control set window size
- connection-oriented:
 - handshaking (exchange of control messages) initializes sender, receiver state before data exchange
- flow controlled:
 - sender will not overwhelm receiver

TCP segment structure



설정은 X, 역할은
가장 파일

→ data 시각화

segment seq #: counting

- bytes of data into bytestream
(not segments!) → bytestream
내에서 축소
비트로 표기
(즉, 확장)
- flow control: # bytes
receiver willing to accept
↳ 1번에 할 수 있는 bytes

Message
data sent by
application into
TCP socket

TCP sequence numbers, ACKs

Sequence numbers: each byte has unique sequence number

- byte stream “number” of first byte in segment’s data

Acknowledgements: ACK

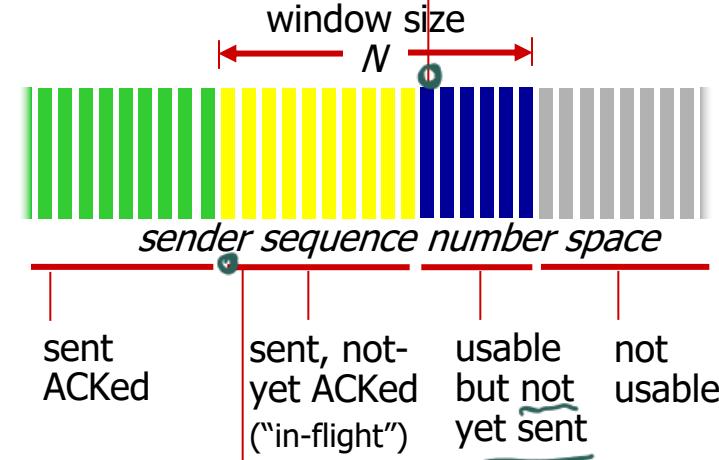
- seq # of next byte expected from other side
- cumulative ACK ↗ 가능 X

Q: how receiver handles out-of-order segments

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

outgoing segment from sender

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

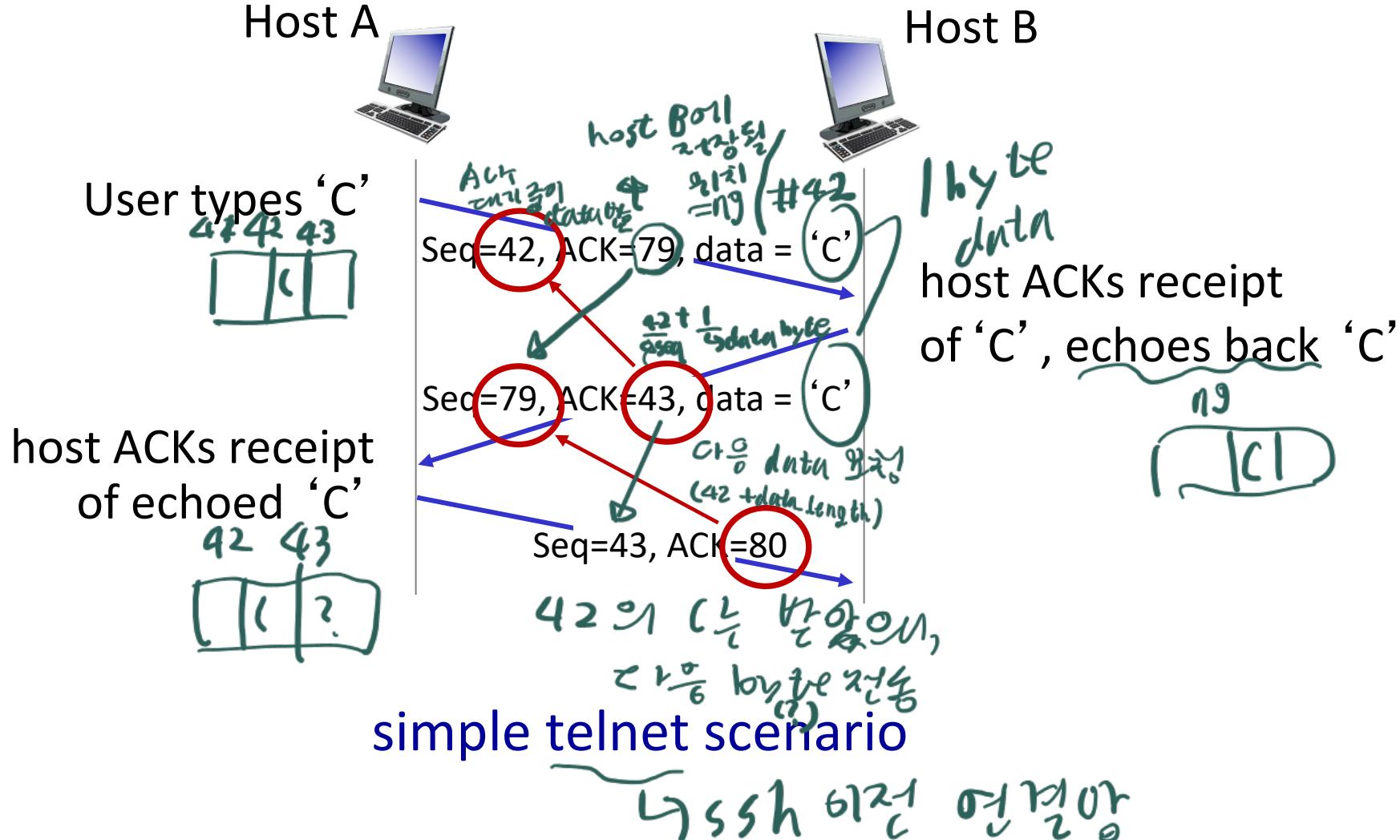


outgoing segment from receiver

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

wrong byte
byte 9/21

TCP sequence numbers, ACKs



TCP round trip time, timeout

TCP timer 길이
too short : multiple retransmit
too long : time loss

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

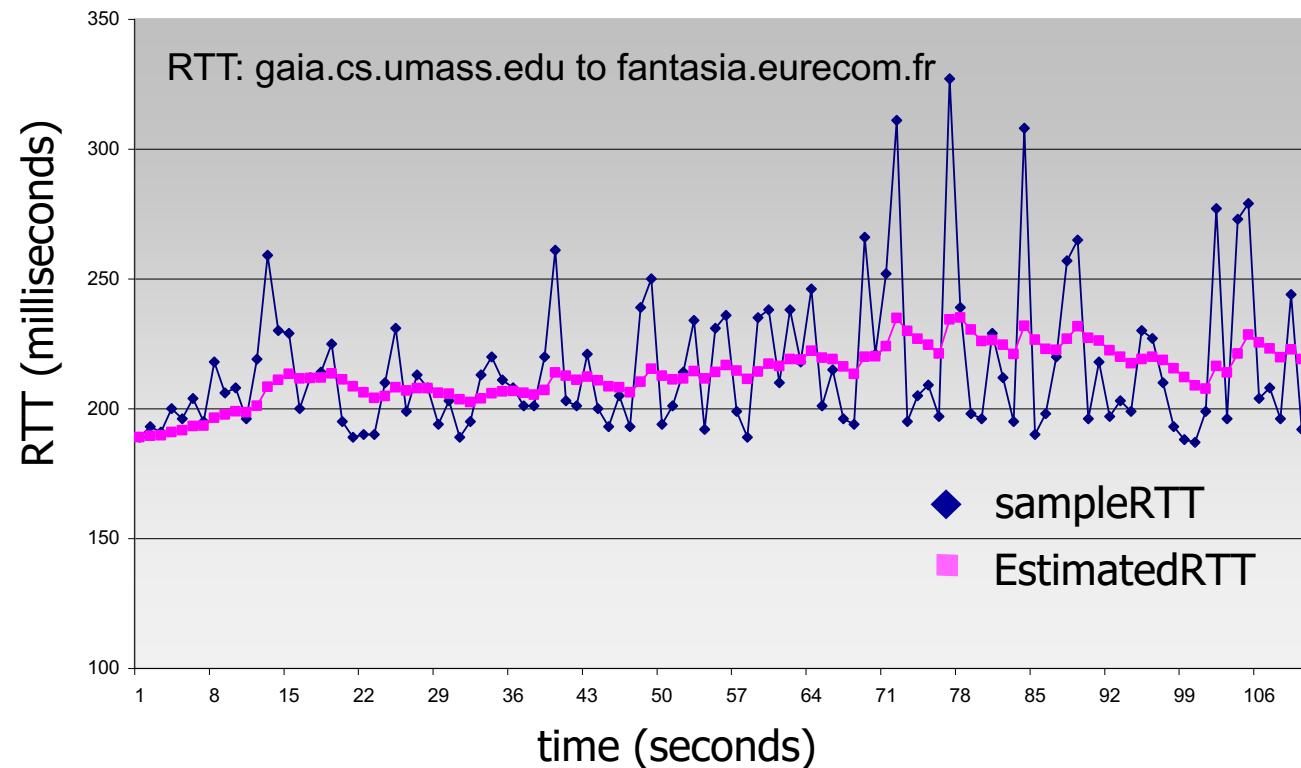
TCP round trip time, timeout

실제 RTT를 예상 RTT에 적용하는 과정
RTT 평균화 과정

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

- exponential weighted moving average (EWMA)
- influence of past sample decreases exponentially fast
- typical value: $\alpha = 0.125$

↑ 시간에 따른
⇒ 가중 평균화
(RTT 평균화)



TCP round trip time, timeout

- timeout interval: **EstimatedRTT** plus “safety margin”
 - large variation in **EstimatedRTT**: want a larger safety margin

$$\text{TimeoutInterval} = \text{EstimatedRTT} + 4 * \text{DevRTT}$$

estimated RTT “safety margin”

- DevRTT: EWMA of **SampleRTT** deviation from **EstimatedRTT**:

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

(typically, $\beta = 0.25$)

이전값과 차이 \Rightarrow 표준차이 \Rightarrow margin↑

TCP Sender (simplified)

event: data received from application

- 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

application

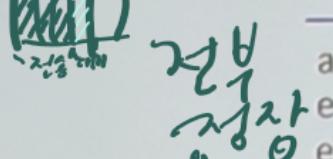
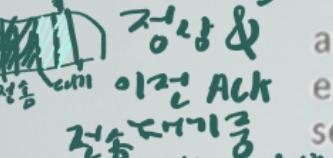
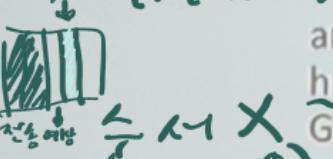
event: timeout

- retransmit segment that caused timeout
- restart timer

event: ACK received

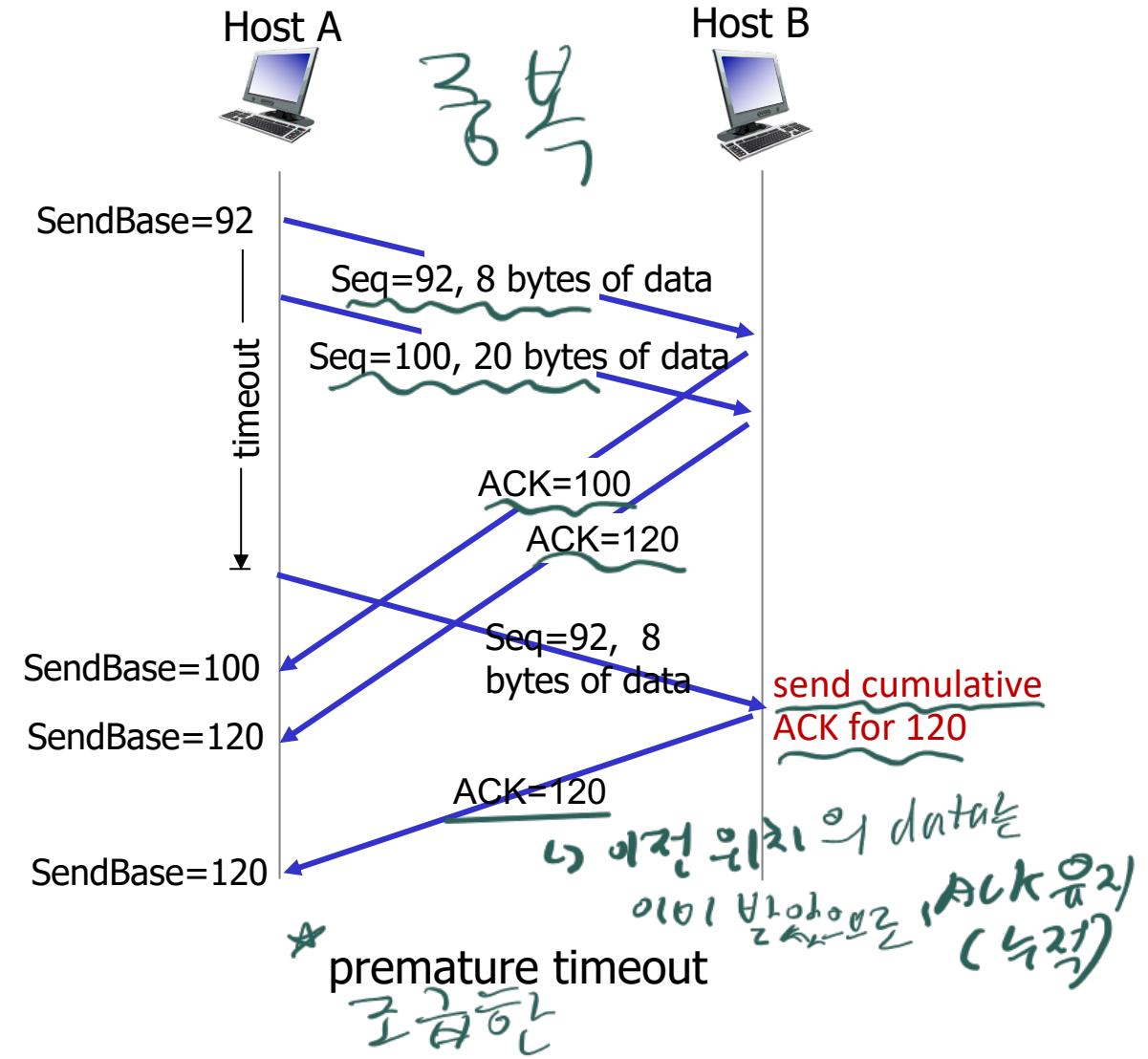
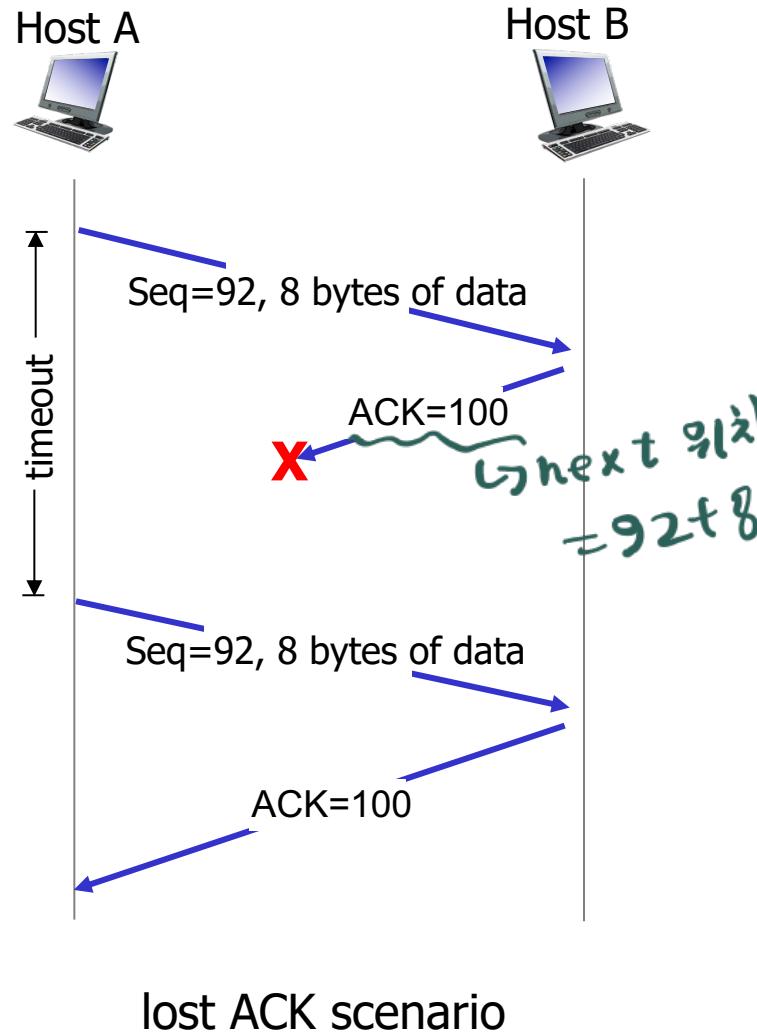
- if ACK acknowledges previously unACKed segments
 - update what is known to be ACKed 성공 update
 - start timer if there are still unACKed segments timer 2ms/2x

TCP Receiver: ACK generation [RFC 5681]

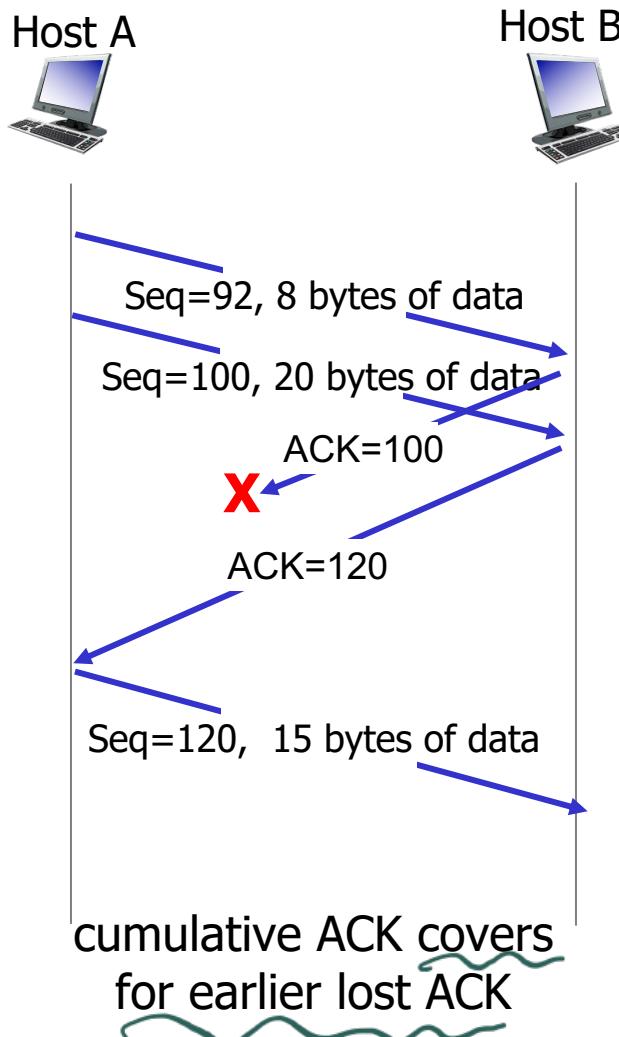
Event at receiver	TCP receiver action
 <p>arrival of <u>in-order</u> segment with expected seq #. All data up to expected seq # already ACKed</p>	<u>delayed ACK</u> . Wait up to 500ms for next segment. If <u>no next segment</u> , send ACK \rightarrow 다음의 packet은 이미 처리된 packet으로 \rightarrow 리소스 (재사용)
 <p>arrival of <u>in-order</u> segment with expected seq #. One other segment has ACK pending</p>	<u>immediately send single cumulative ACK</u> , ACKing both in-order segments
 <p>arrival of <u>out-of-order</u> segment higher-than-expect seq. #. Gap detected</p>	<u>immediately send duplicate ACK</u> , indicating seq. # of next expected byte \rightarrow 허위 ACK
 <p>arrival of segment that partially or completely fills gap</p>	<u>immediate send ACK</u> , provided that segment starts at lower end of gap \rightarrow 실제 ACK



TCP: retransmission scenarios



TCP: retransmission scenarios



누적 ACK를
사용함으로,
이전 Ack까지
그동안 발생한lost
cover 된다

TCP fast retransmit

TCP fast retransmit

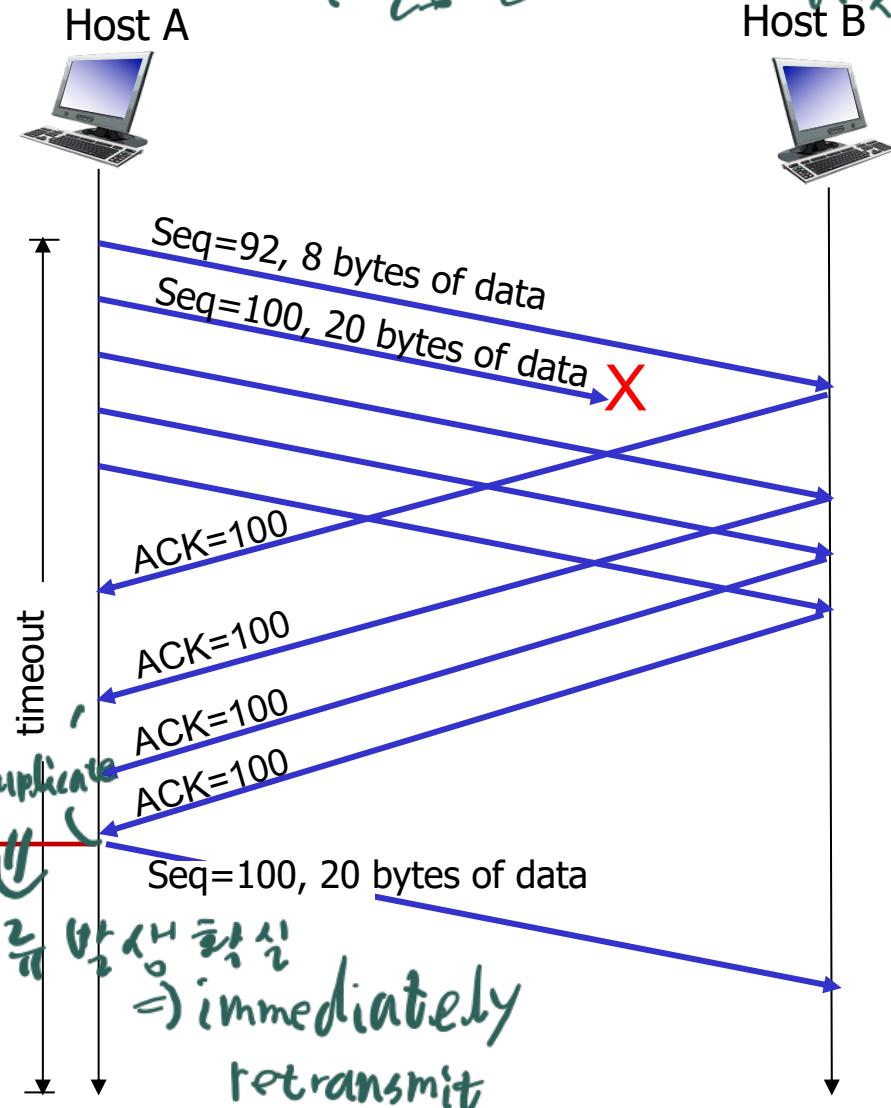
if sender receives 3 additional ACKs for same data (“triple duplicate ACKs”), resend unACKed segment with smallest seq #

- 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!

3 번째 번갈아가며 : so frequently
번갈아가며 : timeout (3 * RTT) 까지 걸립니다
기다리지
기다리지



How to determine TCP window size? ACK 헤더의

경험적 평가?

RTT 평균?

packet loss 허용?

→ 같은 번호의

보내는 두 개의
packet은

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

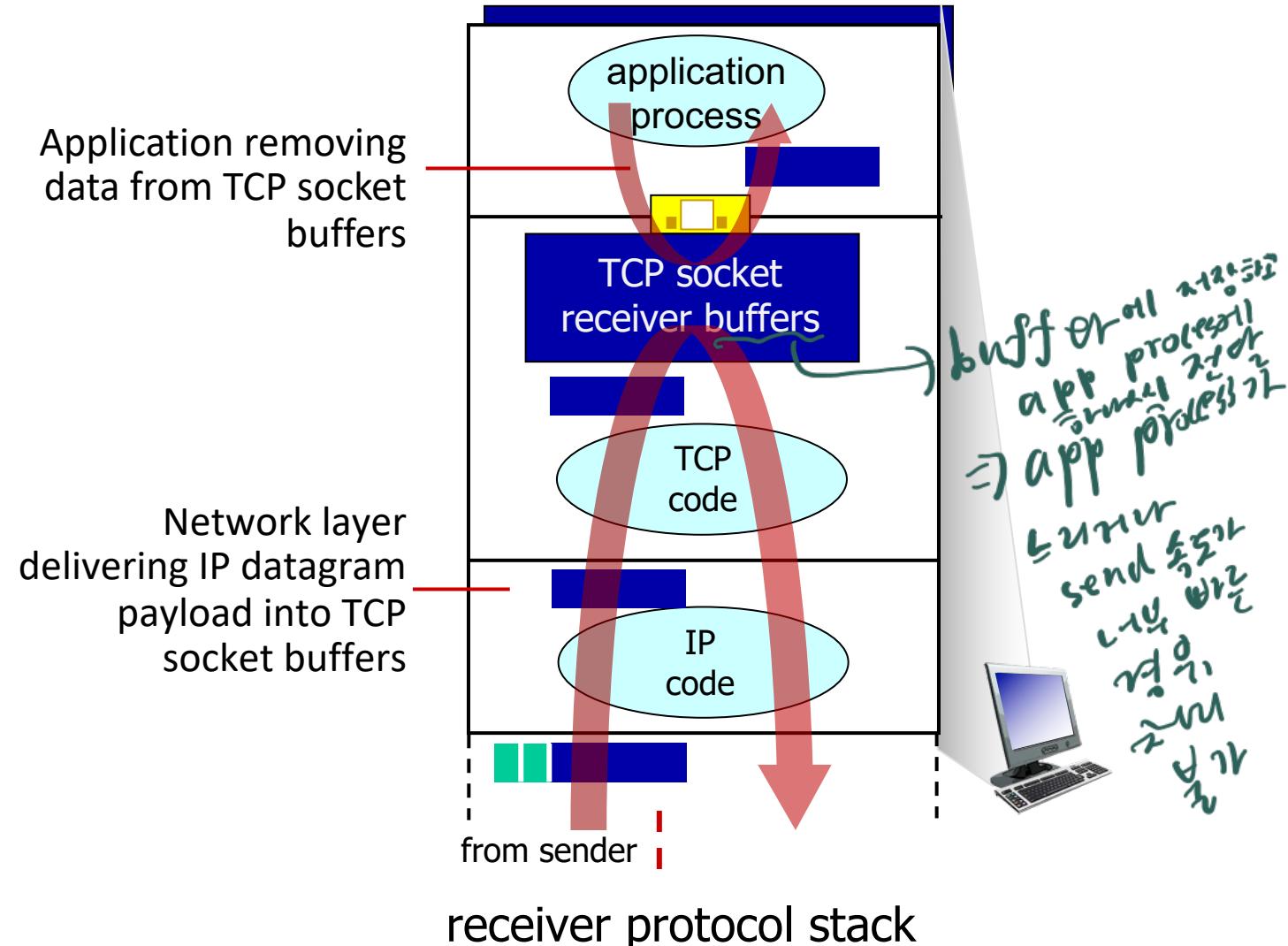


→ OHG 8/21
OVA 2/21
Xplor 2/21

TCP flow control

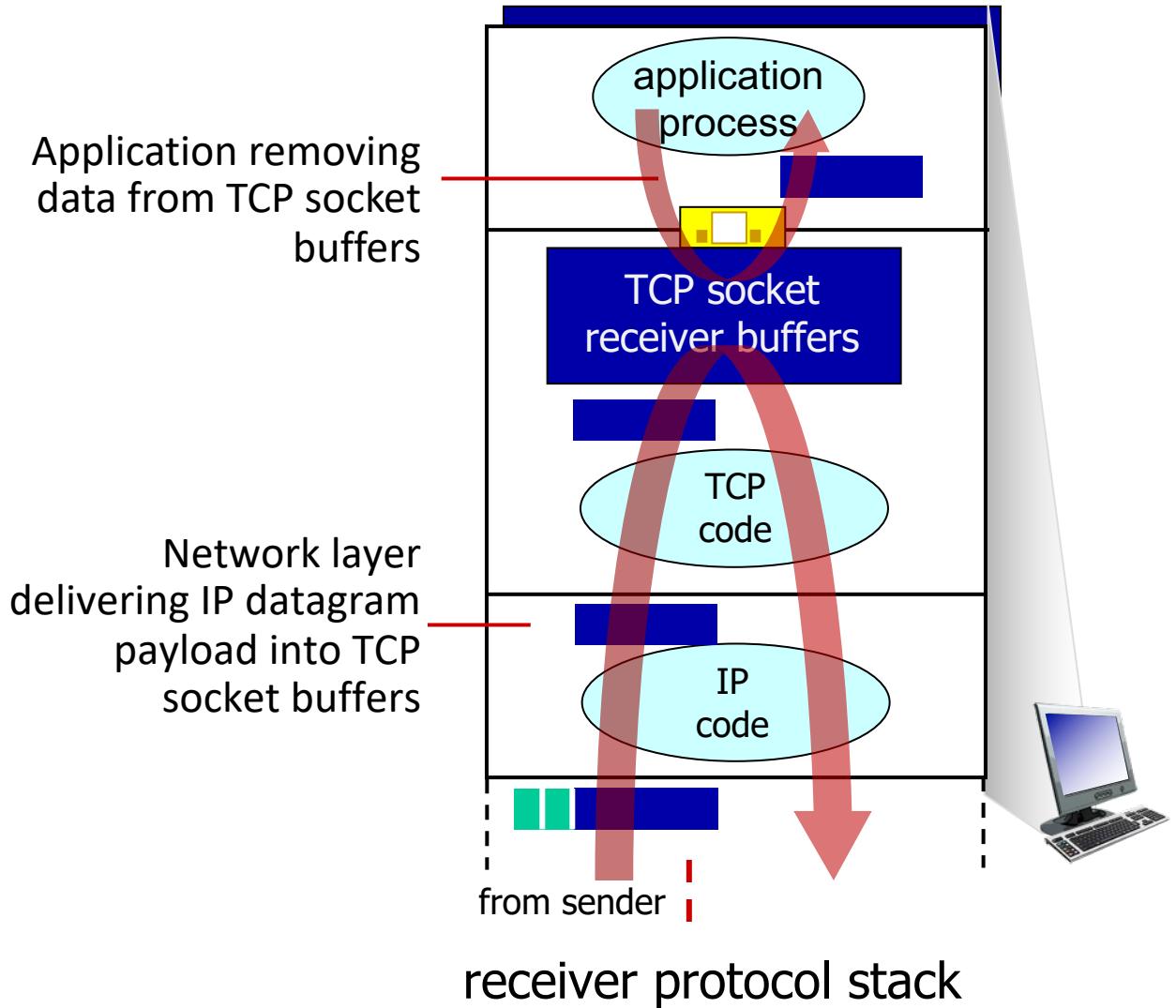
Q: What happens if network layer delivers data faster than application layer removes data from socket buffers?

$\Rightarrow 2^{121} \frac{1}{2} > 2$ (overflow)



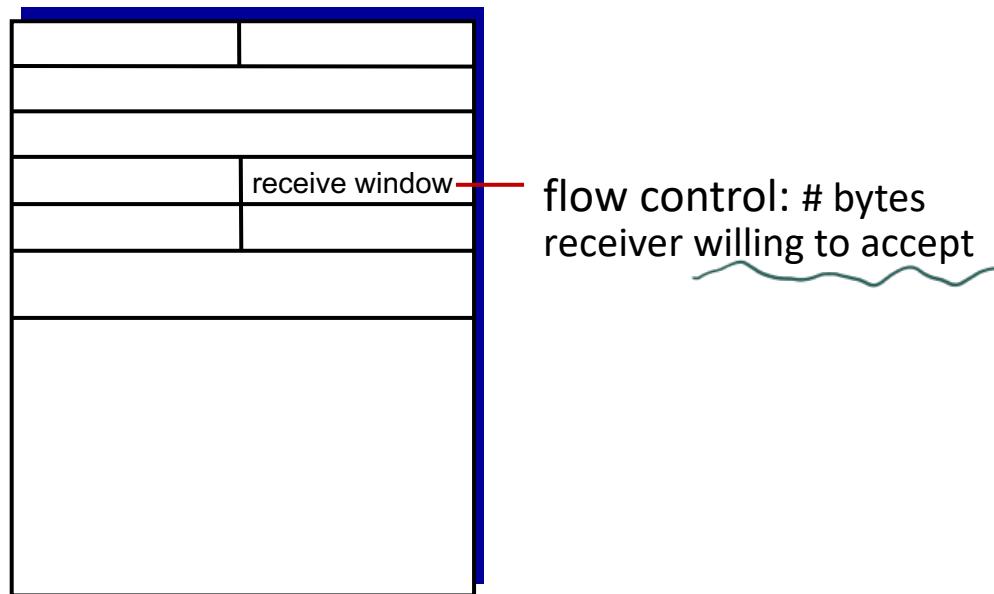
TCP flow control

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

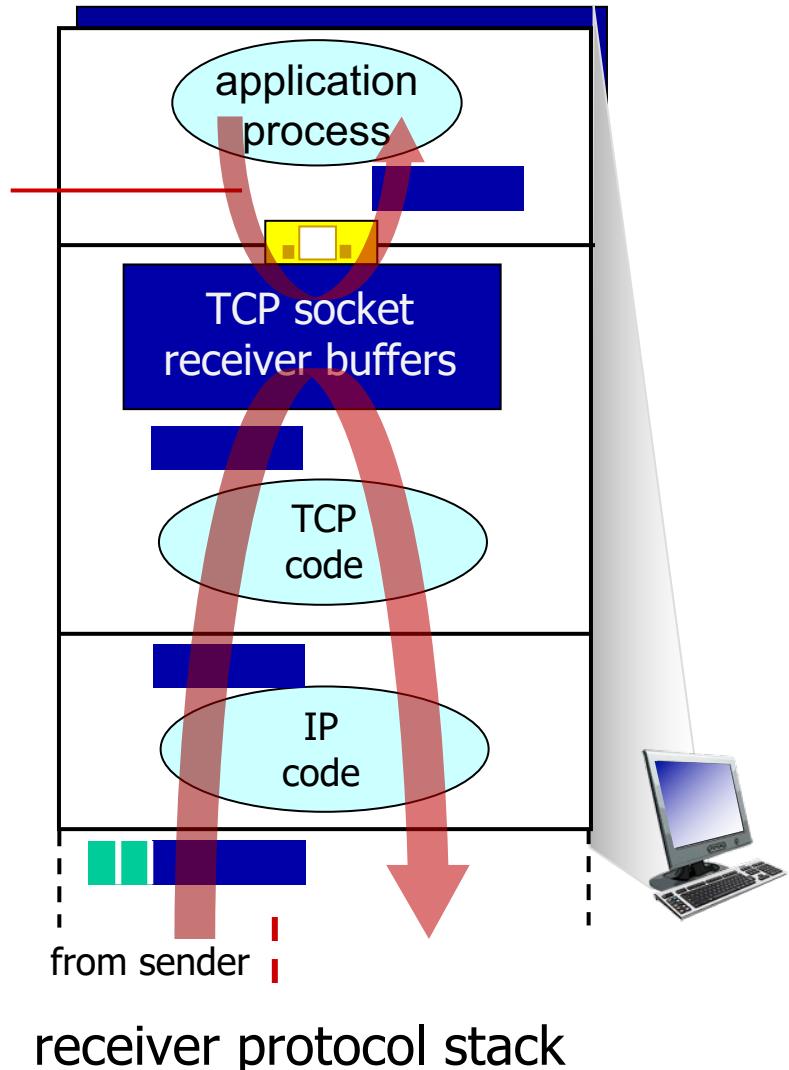


TCP flow control

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



Application removing data from TCP socket buffers

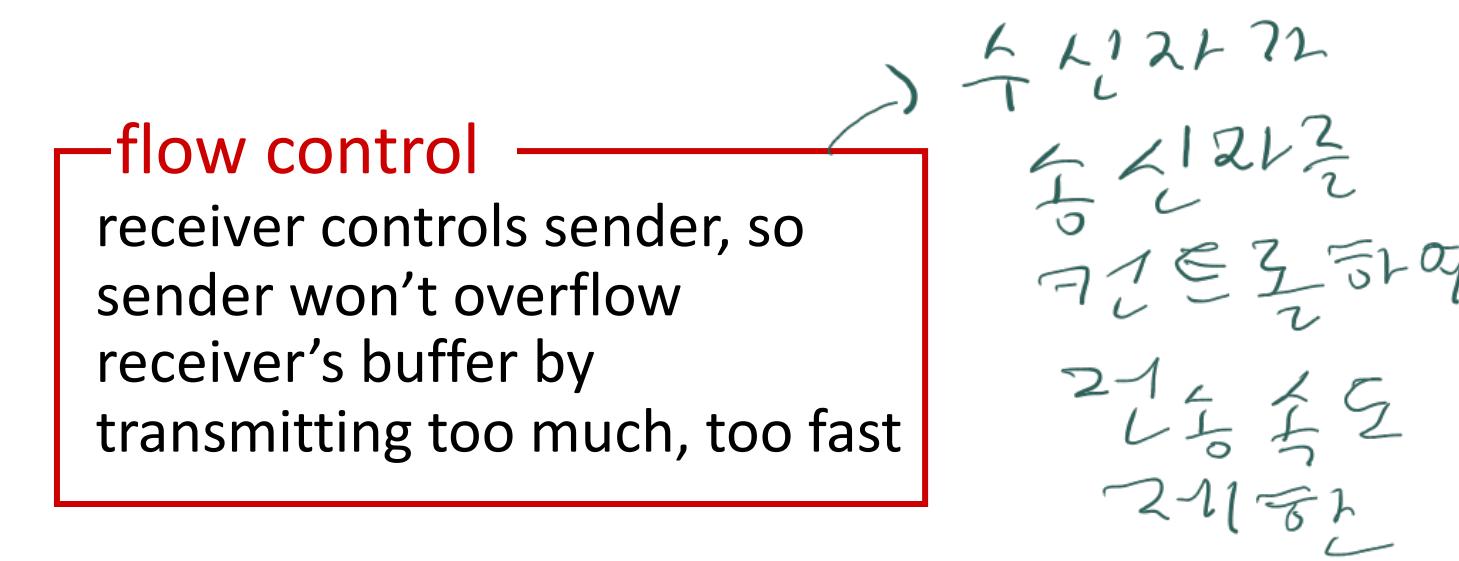


TCP flow control

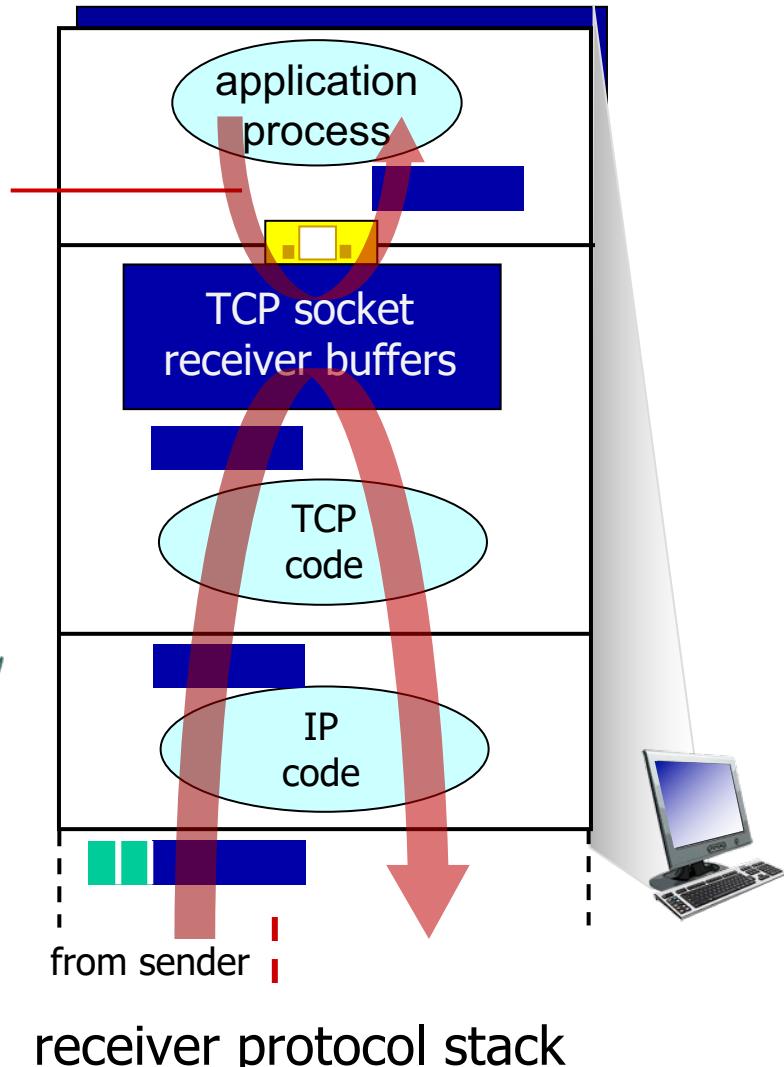
Q: What happens if network layer delivers data faster than application layer removes data from socket buffers?

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



TCP flow control

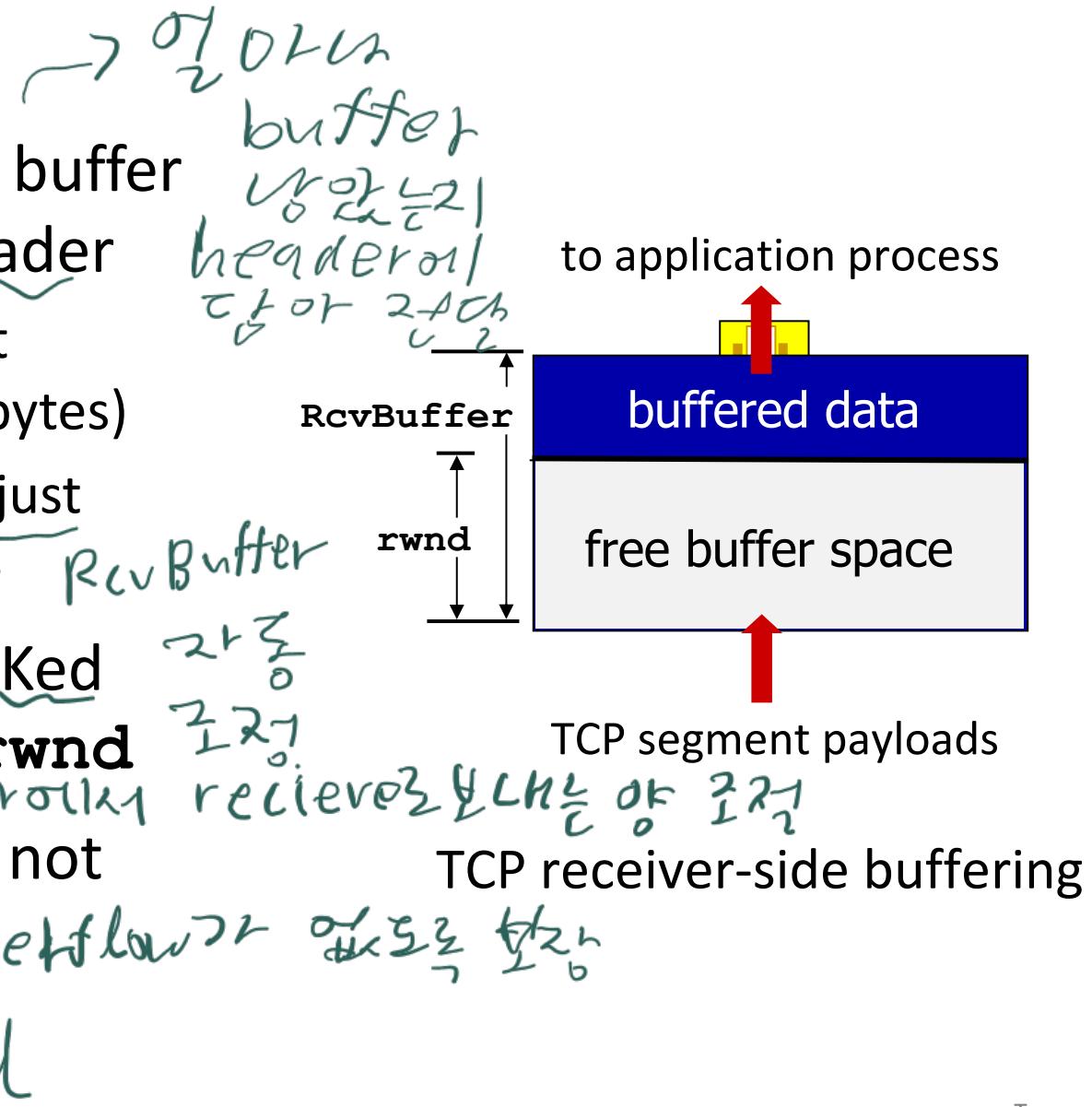
- TCP receiver “advertises” free buffer space in **rwnd** field in TCP header

- **RcvBuffer** size set via socket options (typical default is 4096 bytes)

- many operating systems autoadjust **RcvBuffer**

- sender limits amount of unACKed (“in-flight”) data to received **rwnd**

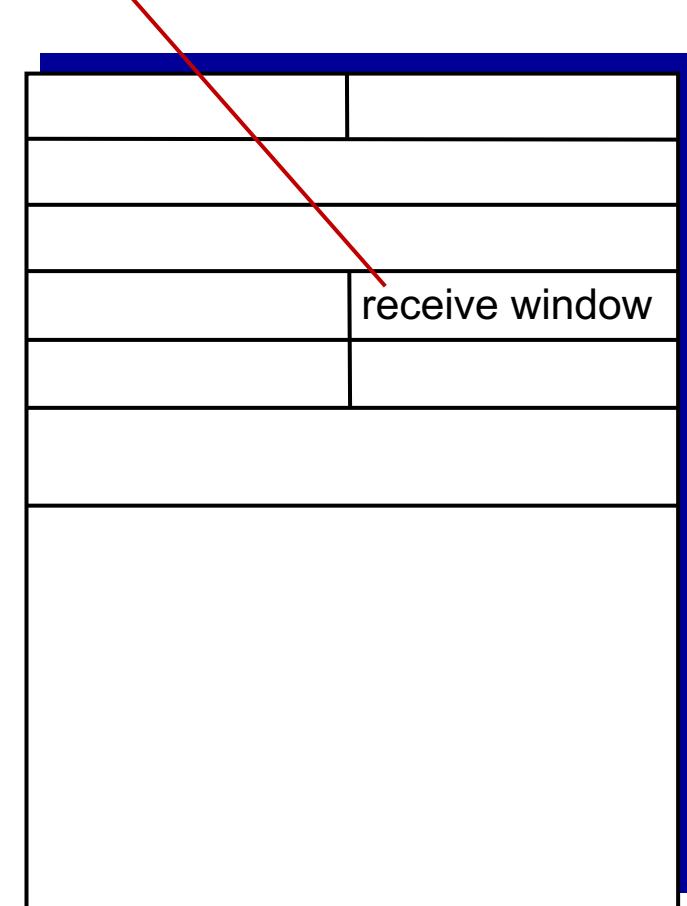
- guarantees receive buffer will not overflow → 수신 버퍼에 overflow가 없도록 보장



TCP flow control

- TCP receiver “advertises” free buffer space in **rwnd** field in TCP header
 - **RcvBuffer** size set via socket options (typical default is 4096 bytes)
 - many operating systems autoadjust **RcvBuffer**
- sender limits amount of unACKed (“in-flight”) data to received **rwnd**
- guarantees receive buffer will not overflow

flow control: # bytes receiver willing to accept



TCP segment format

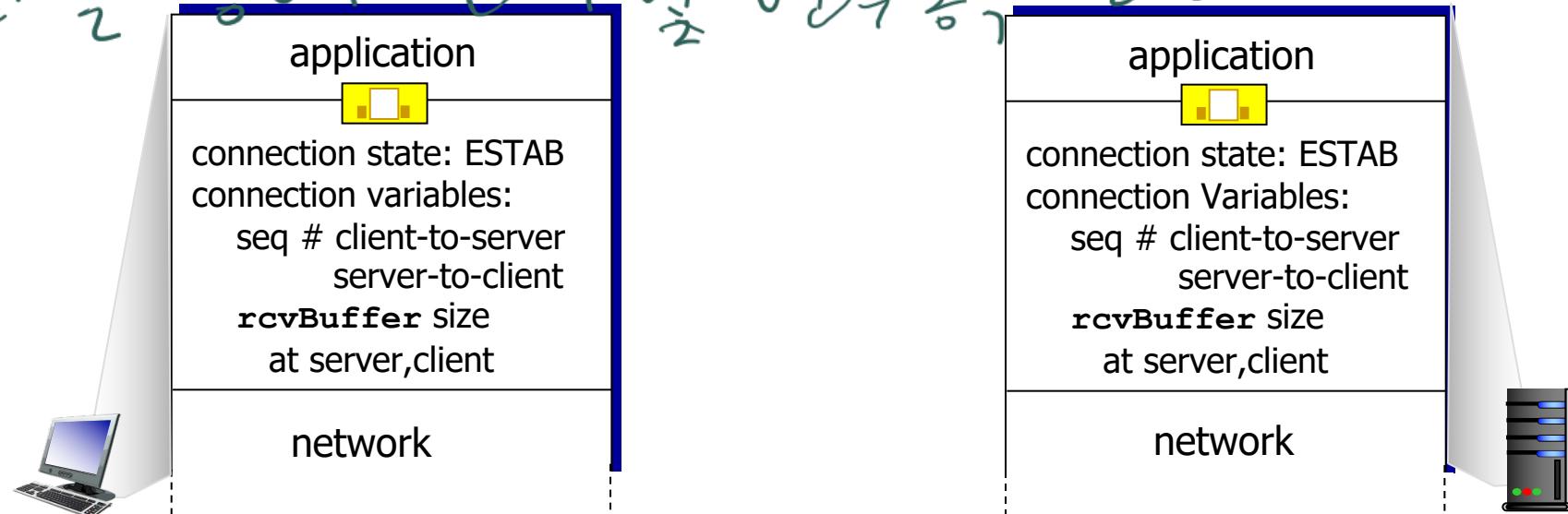
TCP connection management

→ TCP/IP handshake
handshake 필요
(connection 초기화)

before exchanging data, sender/receiver “handshake”:

- agree to establish connection (each knowing the other willing to establish connection)
- agree on connection parameters (e.g., starting seq #s)

⇒ 예전에 상대 인식 이후 서비스 종료 신호



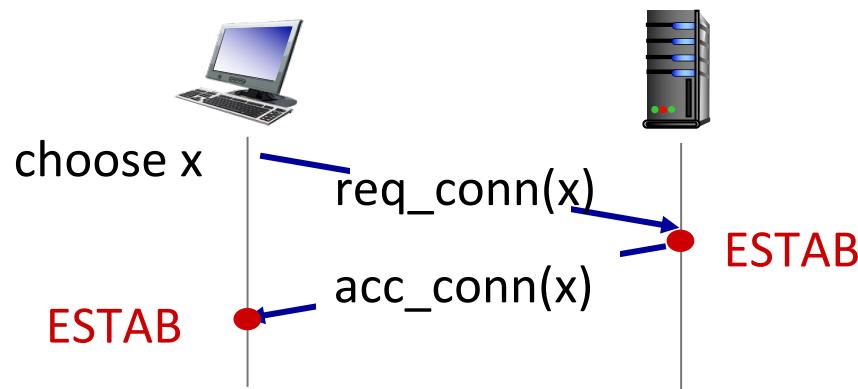
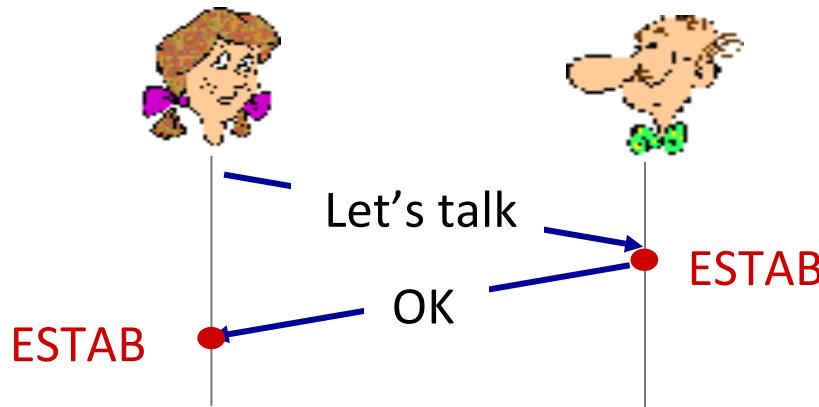
```
Socket clientSocket =  
    newSocket("hostname", "port number");
```

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

Agreeing to establish a connection

On
Lec 10

2-way handshake:



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

- variable delays
- retransmitted messages (e.g. $\text{req_conn}(x)$) due to message loss
- message reordering
- can't "see" other side

TCP 3-way handshake

on
Lec 10
Server state

Client state

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

LISTEN

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

SYNSENT

choose init seq num, x
send TCP SYN msg



SYNbit=1, Seq= x



ESTAB

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

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

ACKbit=1, ACKnum= $y+1$

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

LISTEN

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

SYN RCVD

received ACK(y)
indicates client is live

ESTAB

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

On Lecture 10

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



Principles of congestion control

Congestion:

- informally: “too many sources sending too much data too fast for network to handle”
- manifestations: 
 - long delays (queueing in router buffers)
 - packet loss (buffer overflow at routers)
- different from flow control!
- a top-10 problem!

↳ 중요한 문제 중
② 증후군



congestion control:

too many senders,
sending too fast

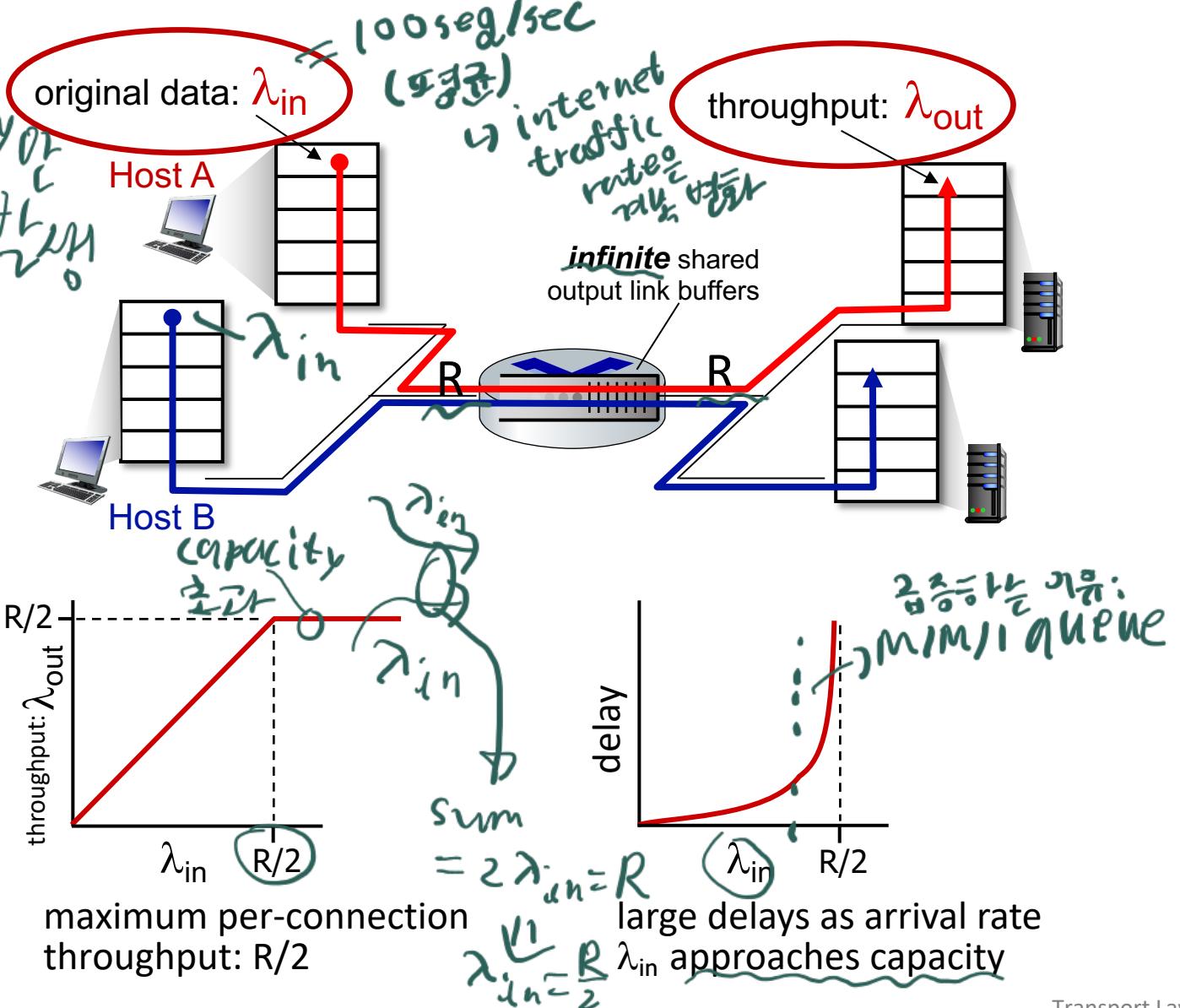
flow control: one sender
too fast for one receiver

Causes/costs of congestion: scenario 1

Simplest scenario:

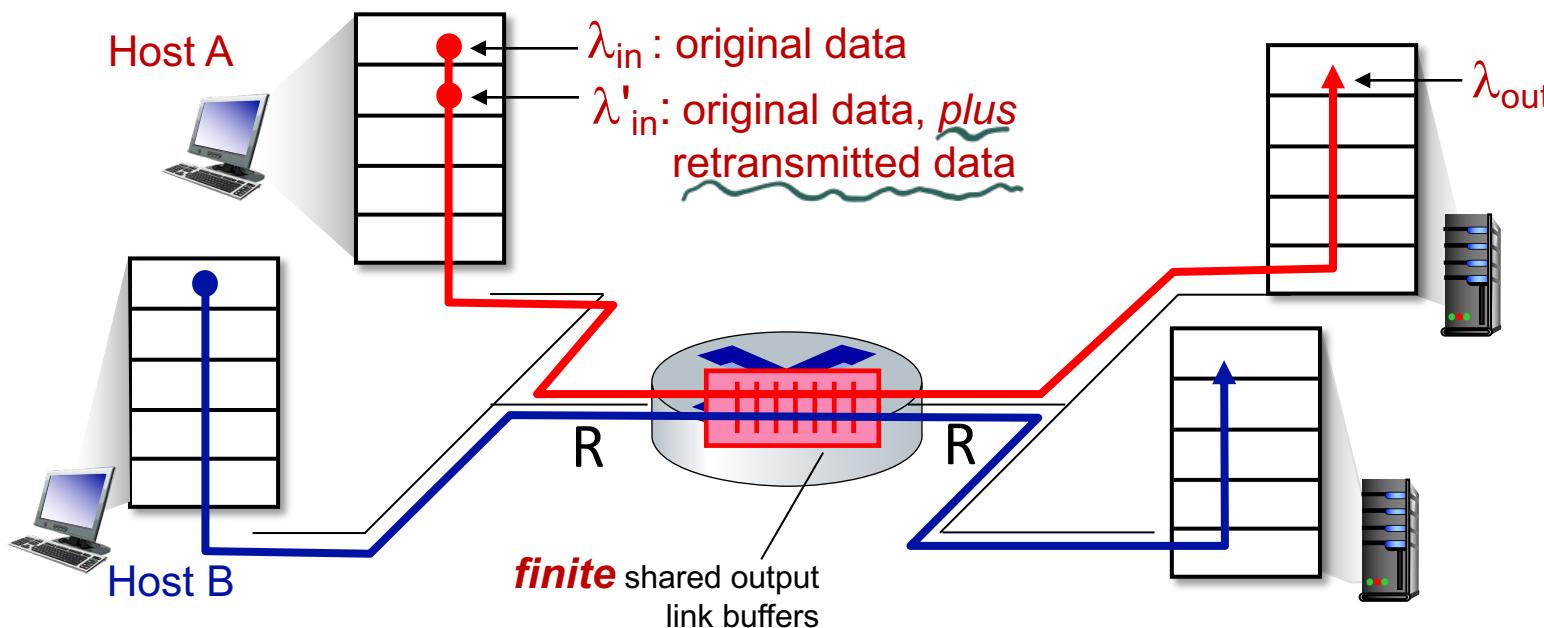
- one router, infinite buffers
- input, output link capacity: R
- two flows
- no retransmissions needed

Q: What happens as arrival rate λ_{in} approaches $R/2$?

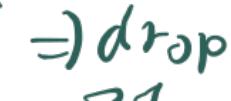


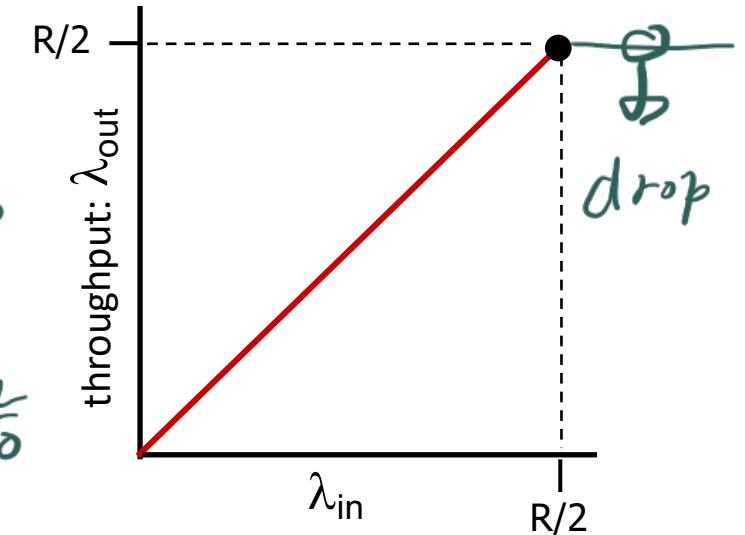
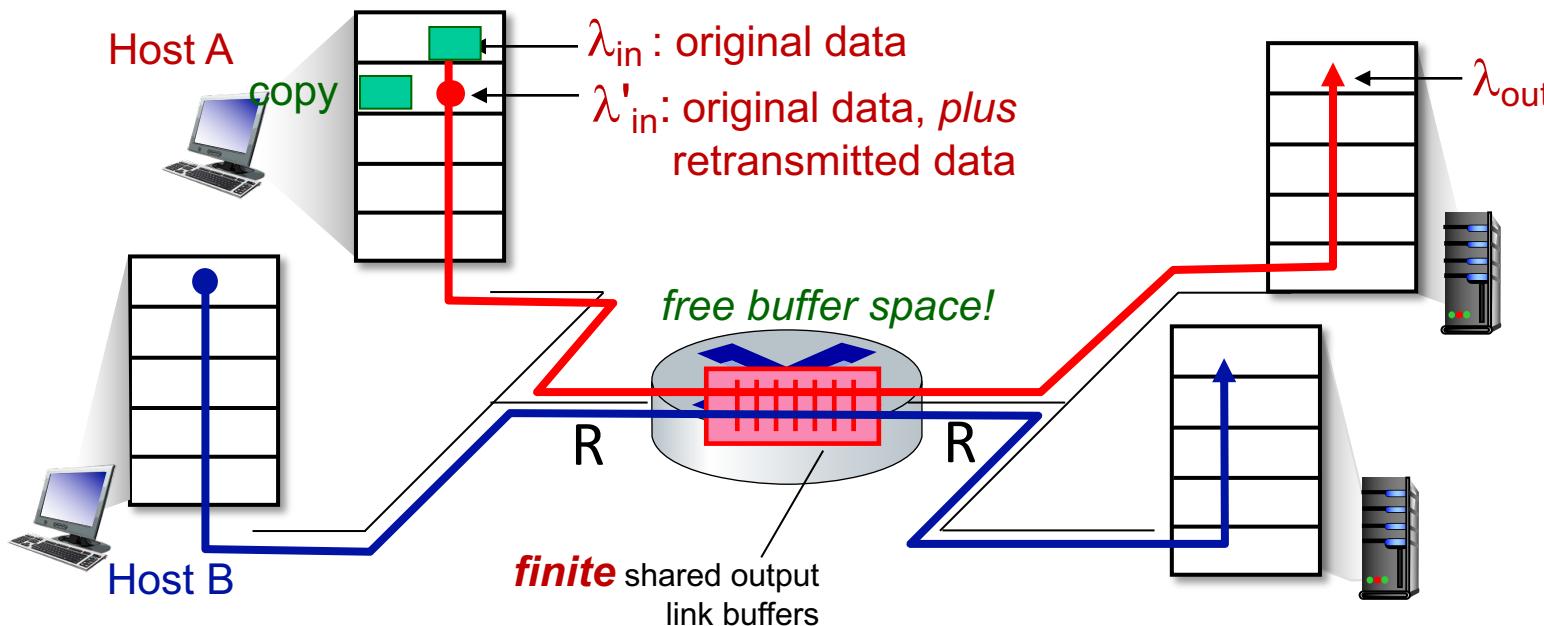
Causes/costs of congestion: scenario 2

- one router, finite buffers *-lost ΗLηυ*
- sender retransmits lost, timed-out packet
 - application-layer input = application-layer output: $\lambda_{in} = \lambda_{out}$
 - transport-layer input includes *retransmissions* : $\lambda'_{in} \geq \lambda_{in}$



Causes/costs of congestion: scenario 2

Idealization: perfect knowledge → loss 
→ drop 
→ drop 
→ drop 
→ drop 



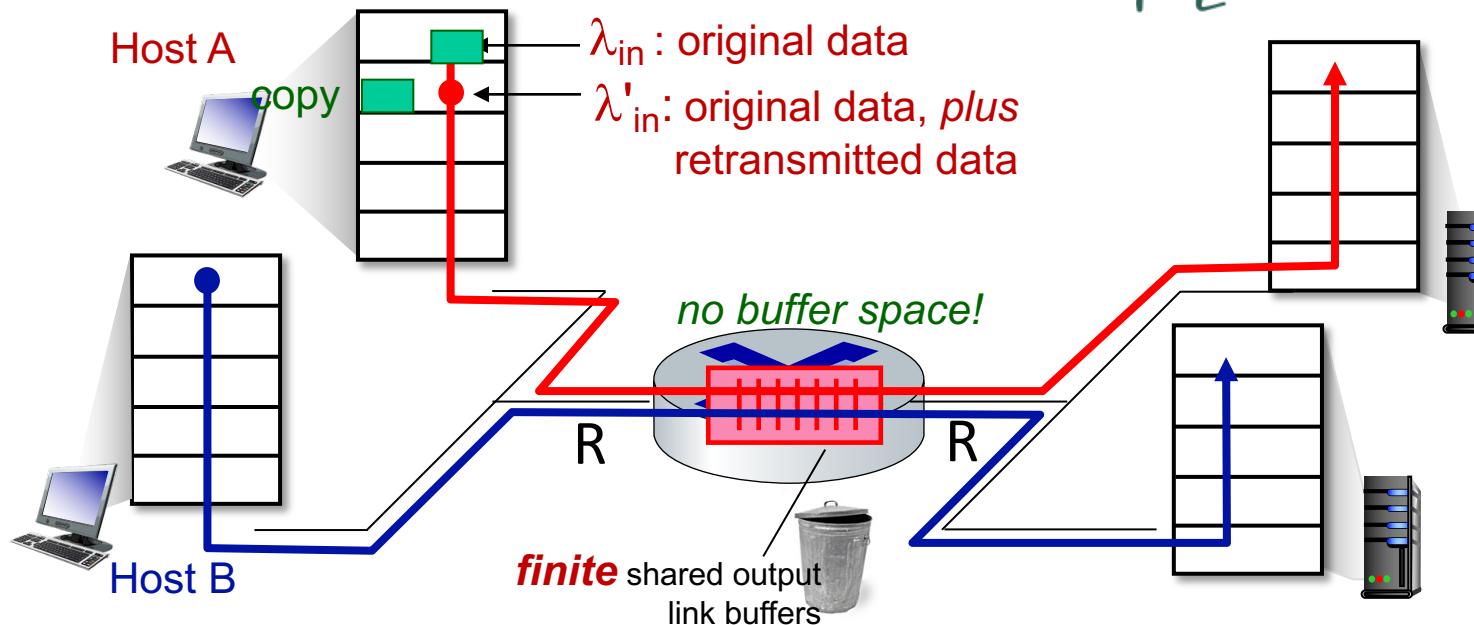
Causes/costs of congestion: scenario 2

Idealization: *some* perfect knowledge

- packets can be lost (dropped at router) due to full buffers
- sender knows when packet has been dropped: only resends if packet known to be lost

→ loss 발생 예측
2 번 충돌
⇒ lost 발생 가능

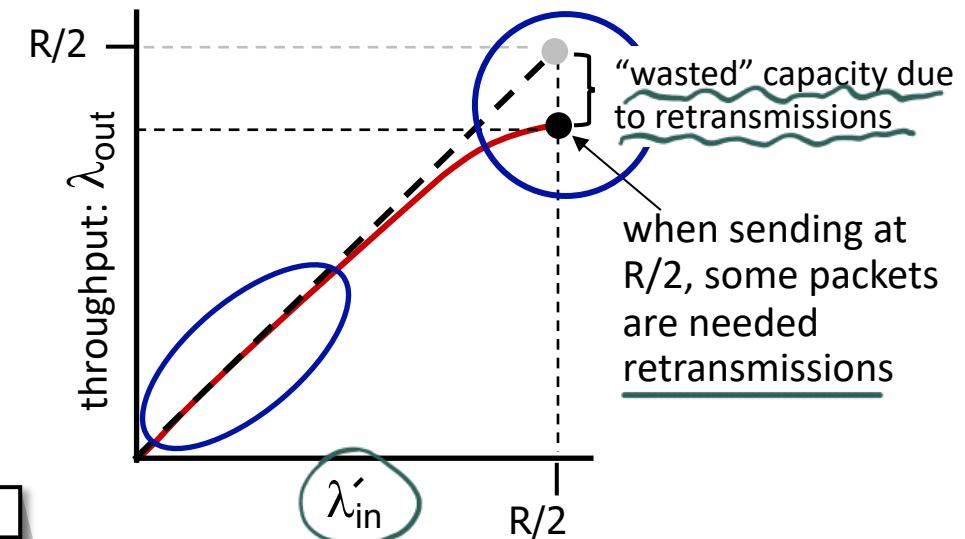
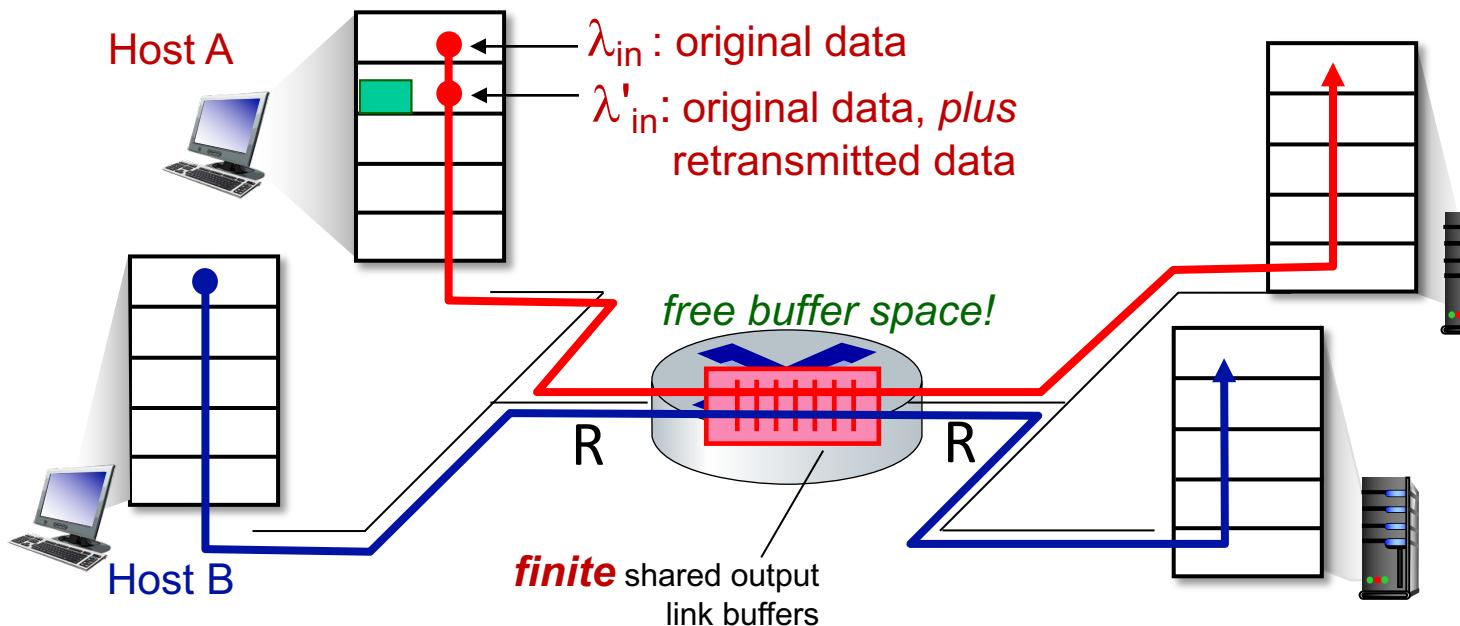
↳ drop된 경우 인지 가능 ⇒ 재전송



Causes/costs of congestion: scenario 2

Idealization: *some* perfect knowledge

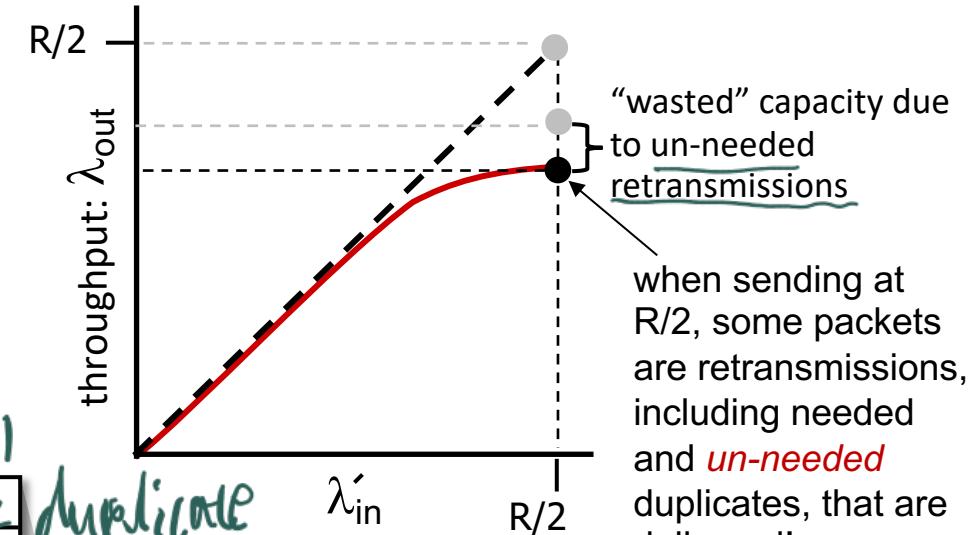
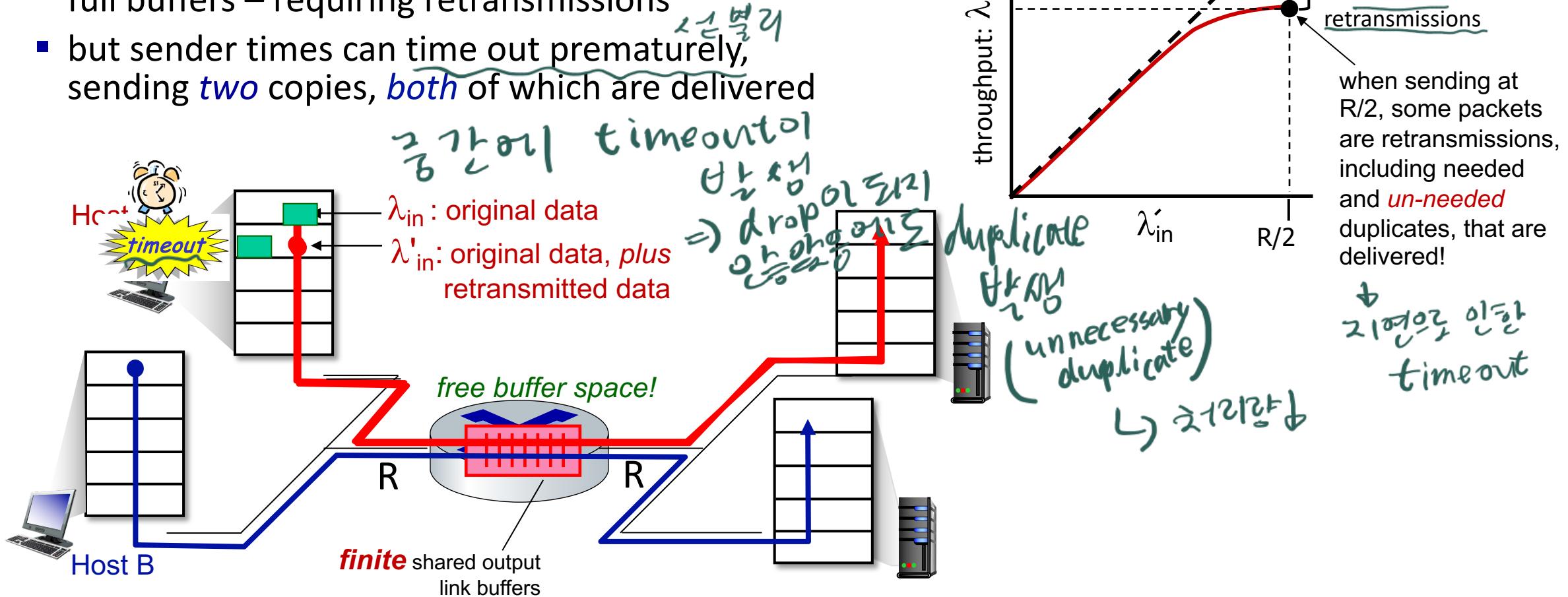
- packets can be lost (dropped at router) due to full buffers
- sender knows when packet has been dropped: only resends if packet *known* to be lost



Causes/costs of congestion: scenario 2

Realistic scenario: un-needed duplicates

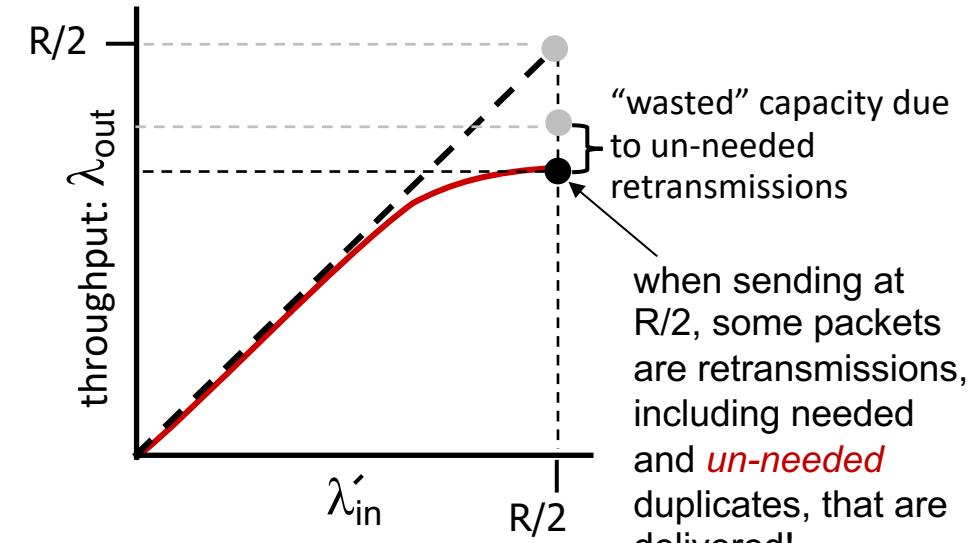
- packets can be lost, dropped at router due to full buffers – requiring retransmissions
- but sender times can time out prematurely, sending **two** copies, **both** of which are delivered



Causes/costs of congestion: scenario 2

Realistic scenario: *un-needed duplicates*

- packets can be lost, dropped at router due to full buffers – requiring retransmissions
- but sender times can time out prematurely, sending *two* copies, *both* of which are delivered



"costs" of congestion:

- more work (retransmission) for given receiver throughput
- unneeded retransmissions: link carries multiple copies of a packet
 - decreasing maximum achievable throughput

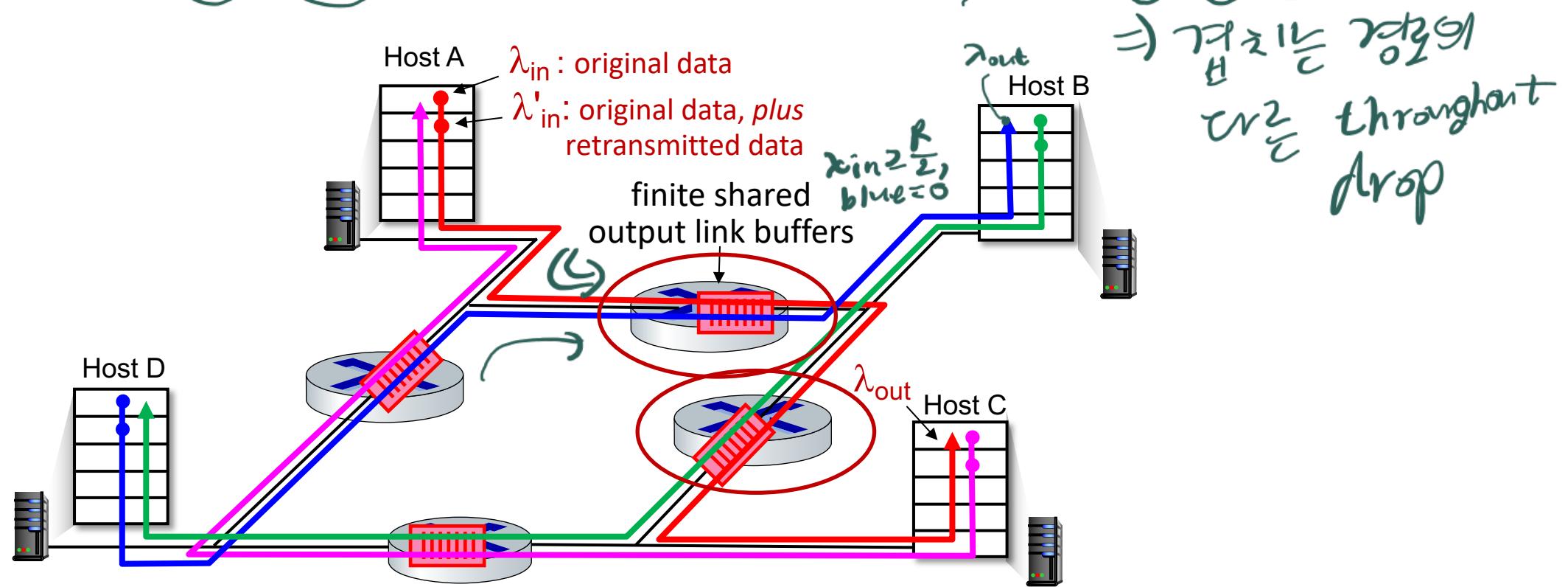
↳ 부적절한 재전송은 최대가능률을 ↓

Causes/costs of congestion: scenario 3

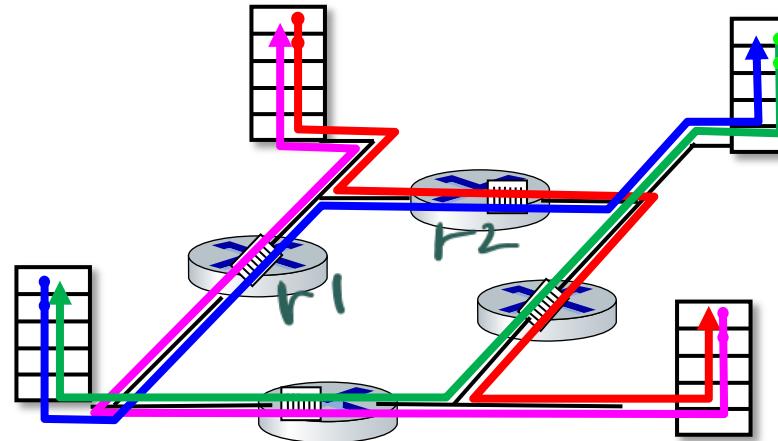
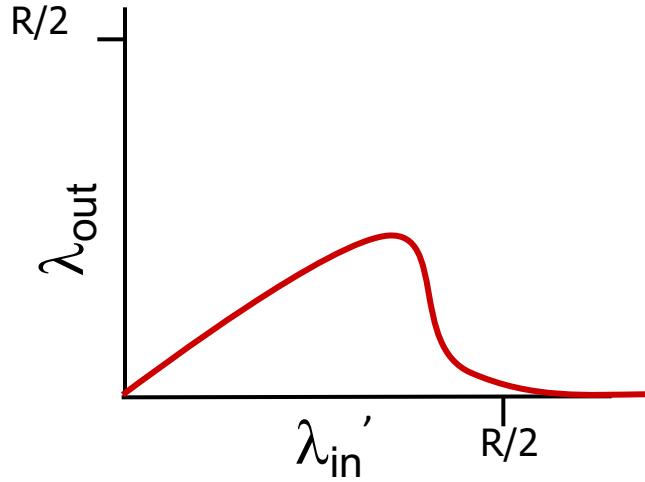
- four senders
- multi-hop paths
- timeout/retransmit

Q: what happens as λ_{in} and λ'_{in} increase ?

A: as red λ'_{in} increases, all arriving blue pkts at upper queue are dropped, blue throughput $\rightarrow 0$



Causes/costs of congestion: scenario 3



blue 는 r1% 차원으로,
but r2에서-에서 drop된다
r1 차원도 낭비

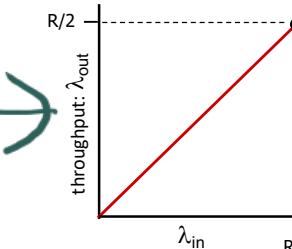
another “cost” of congestion:

- when packet dropped, any upstream transmission capacity and buffering used for that packet was wasted!

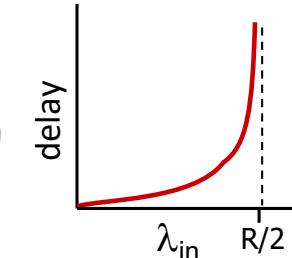
Causes/costs of congestion: insights

- throughput can never exceed capacity

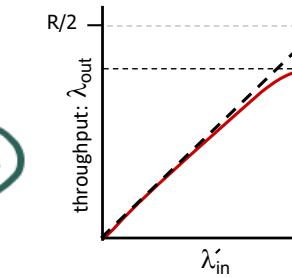
⇒ $\text{전송량} \leq \text{용량}$



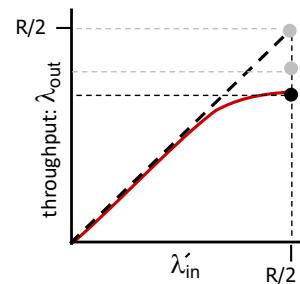
- delay increases as capacity approached



- loss/retransmission decreases effective throughput

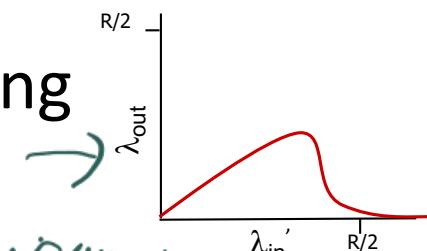


- un-needed duplicates further decreases effective throughput



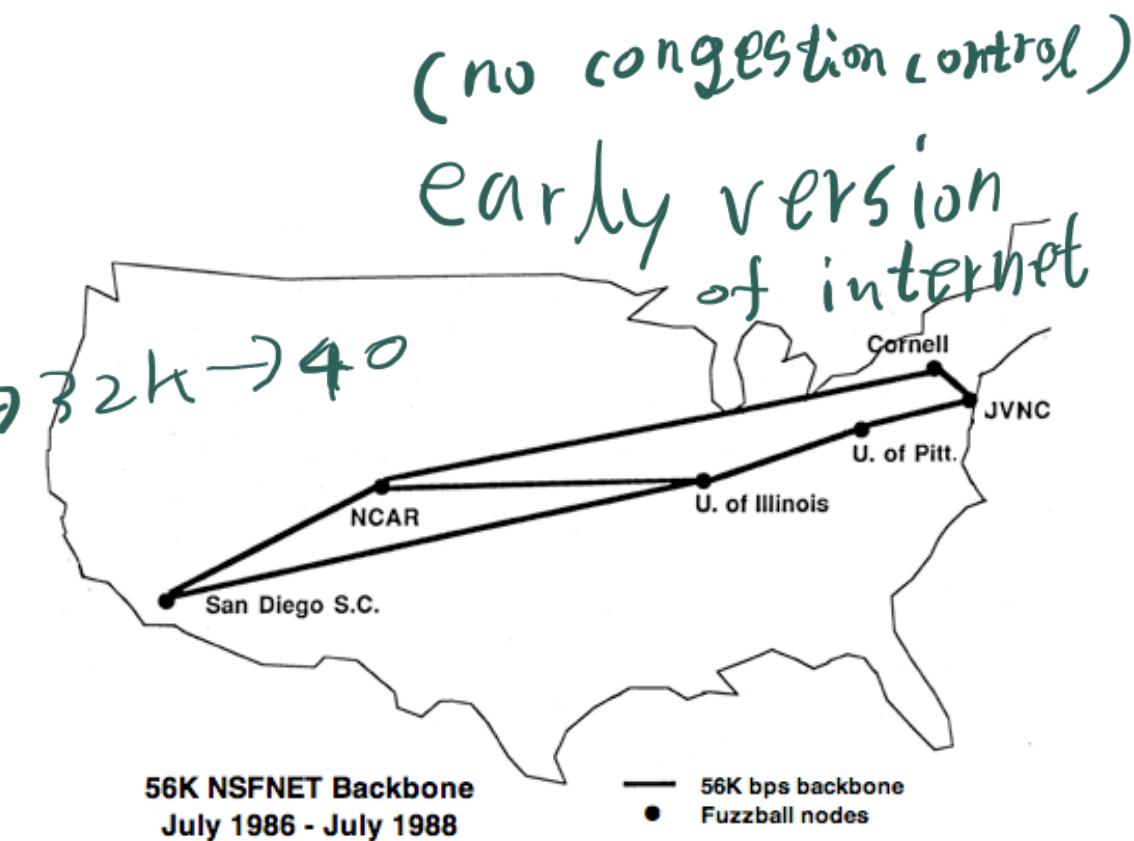
- upstream transmission capacity / buffering wasted for packets lost downstream

↳ downstream packets lost \Rightarrow upstream packets wasted



“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



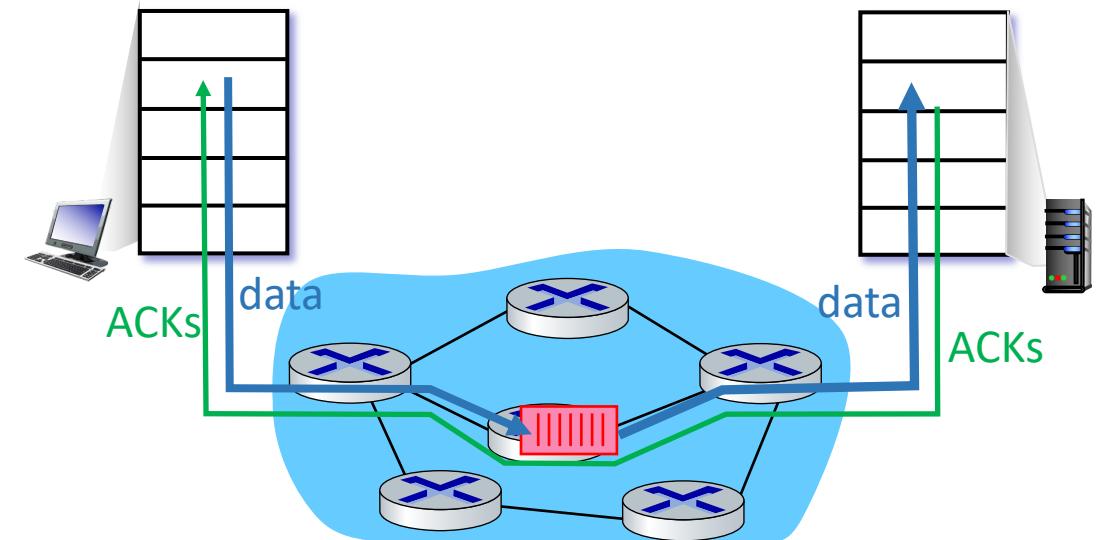
Approaches towards congestion control

End-end congestion control:

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

⇒ congestion *at the end hosts*

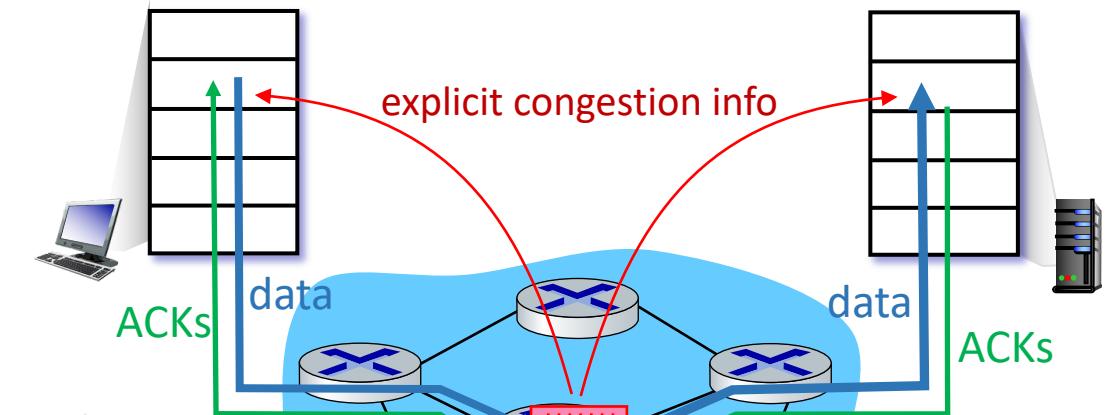
of the end hosts loss, delay *is estimated*



Approaches towards congestion control

Network-assisted congestion control:

- routers provide direct feedback to sending/receiving hosts with flows passing through congested router
2141822 → congestion level / sending rate 설정 가능
- may indicate congestion level or explicitly set sending rate
2141822
- TCP ECN, ATM, DECbit protocols



Next...

- *Chapter 3.7 TCP Congestion Control*