

Advanced Networks

Transport layer: TCP

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Dr. Wei Bao | Lecturer
School of Computer Science



THE UNIVERSITY OF
SYDNEY

Go-back-N:

- › sender can have up to N unacked packets in pipeline
- › receiver only sends **cumulative ack**
 - does not ack packet if there is a gap
- › sender has timer for oldest unacked packet
 - when timer expires, retransmit *all* unacked packets

Selective Repeat:

- › sender can have up to N unacked packets in pipeline
- › receiver sends **individual ack** packet
- › Add WeChat `edu_assist_pro` contains timer for unacked packet
 - when timer expires, retransmit only that unacked packet



- › “window” of up to N, consecutive unacked pkts allowed

Assignment Project Exam Help

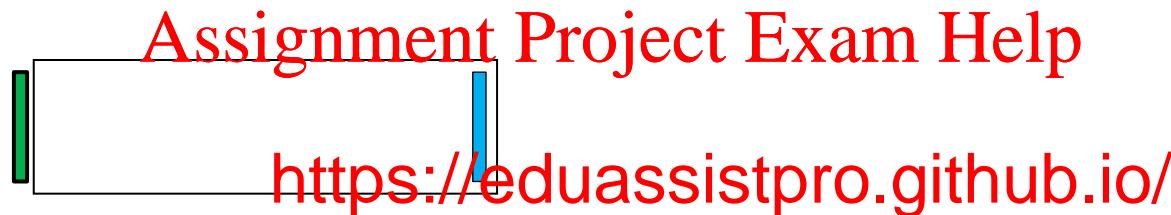
<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

- ❖ ACK(n):ACKs all pkts up to, including seq # n - “*cumulative ACK*”
 - may receive duplicate ACKs (see receiver)
- ❖ timer for oldest in-flight pkt
- ❖ *timeout(n)*: retransmit packet n and all higher seq # pkts in window



- › “window” of up to N, consecutive unacked pkts allowed



Add WeChat `edu_assist_pro`

- ❖ ACK(n): ACKs all pkts up to, including seq # n - “*cumulative ACK*”
 - may receive duplicate ACKs (see receiver)
- ❖ timer for oldest in-flight pkt
- ❖ $\text{timeout}(n)$: retransmit packet n and all higher seq # pkts in window

```
rdt_send(data)  
if (nextseqnum < base+N) {  
    sndpkt[nextseqnum] = make_pkt(nextseqnum,data,chksum)  
    udt_send(sndpkt[nextseqnum])  
    if (base == nextseqnum)
```

Assignment Project Exam Help

Λ
base=1
nextseqnum=1

`rdt_rcv(rcvpkt)`
`&& corrupt(rcvpkt)`

A diagram illustrating a state transition. It features a large circle labeled "Wait" at the top right. An arrow points from the bottom left towards this circle. A large red arrow points from the top left towards the word "Add WeChat".

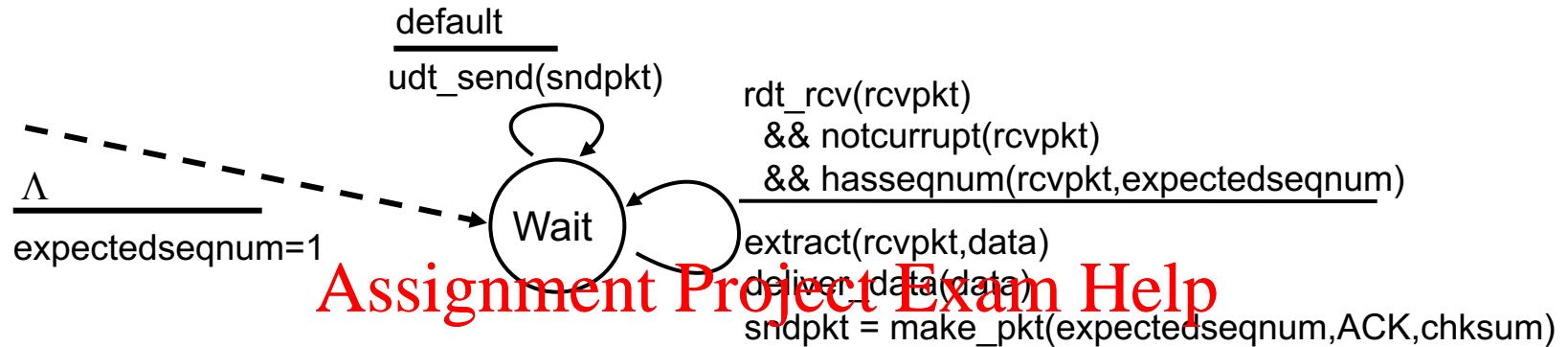
`rdt_rcv(rcvpkt) &&
notcorrupt(rcvpkt)`

```
base = getacknum(rcvpkt)+1  
If (base == nextseqnum)  
    stop_timer  
else  
    start_timer
```

<https://eduassistpro.github.io/>

**St_{base})
[base+1])**

```
...  
udt_send(sndpkt[nextseqnum-1])
```



Assignment Project Exam Help

<https://eduassistpro.github.io/>

ACK-only: always send ACK
highest *in-order* seq #

ly-received pkt with

- may generate duplicate ACKs
 - need only remember **expectedseqnum**
- › out-of-order pkt:
- discard (don't buffer): *no receiver buffering!*
 - re-ACK pkt with highest in-order seq #

sender window (N=4)

0	1	2	3	4	5	6	7	8
0	1	2	3	4	5	6	7	8
0	1	2	3	4	5	6	7	8
0	1	2	3	4	5	6	7	8

sender

send pkt0
send pkt1
send pkt2
send pkt3

(wait)

receiver

receive pkt0, send ack0
receive pkt1, send ack1

Help

receive pkt3, discard,
(re)send ack1

<https://eduassistpro.github.io/>

ceive pkt4, discard,
(re)send ack1
ceive pkt5, discard,
(re)send ack1

0	1	2	3	4	5	6	7	8
0	1	2	3	4	5	6	7	8

rcv ack0, s
rcv ack1, s

ignore duplicate ACK



pkt 2 timeout

0	1	2	3	4	5	6	7	8
0	1	2	3	4	5	6	7	8
0	1	2	3	4	5	6	7	8
0	1	2	3	4	5	6	7	8

send pkt2
send pkt3
send pkt4
send pkt5

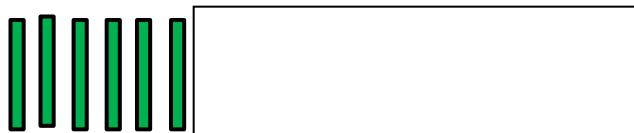
rcv pkt2, deliver, send ack2
rcv pkt3, deliver, send ack3
rcv pkt4, deliver, send ack4
rcv pkt5, deliver, send ack5



- › receiver *individually* acknowledges all correctly received pkts
 - buffers pkts as needed, for eventual in-order delivery to upper layer
- › sender only re <https://eduassistpro.github.io/> received
 - sender timer for each unACKed
- › sender window
- › receiver window



Selective repeat: sender, receiver windows



Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro



sender**data from above:**

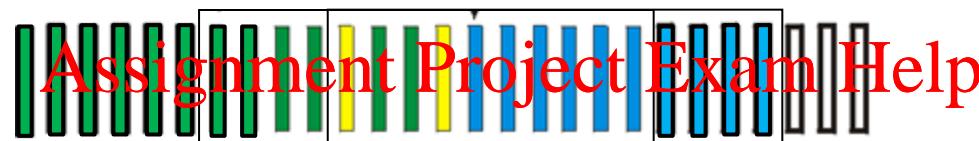
- › if next available seq # in window, send pkt
- time out(n):
 - › resend pkt n, res
 - ACK(n) in [sendbase,sendbase+N-1]:
 - › mark pkt n as received
 - › if n is smallest unACKed pkt, advance window base to next unACKed seq #

receiver**pkt n in [rcvbase, rcvbase+N-1]**

- ❖ send ACK(n)
- out-of-order: buffer
 - der: deliver (also buffered, in-order window to t-yet-received pkt [rcvbase-N,rcvbase-1])
- ACK(n)
- otherwise:
 - ❖ ignore

<https://eduassistpro.github.io/>

Add WeChat **edu_assist_pro**



<https://eduassistpro.github.io/>



Selective repeat in action

sender window (N=4)

0	1	2	3	4	5	6	7	8
0	1	2	3	4	5	6	7	8
0	1	2	3	4	5	6	7	8
0	1	2	3	4	5	6	7	8

sender

send pkt0
send pkt1
send pkt2
send pkt3

(wait)

receiver

receive pkt0, send ack0
receive pkt1, send ack1

receive pkt3, buffer,
send ack3

receive pkt4, buffer,
send ack4
receive pkt5, buffer,
send ack5

0	1	2	3	4	5	6	7	8
0	1	2	3	4	5	6	7	8

0	1	2	3	4	5	6	7	8
---	---	---	---	---	---	---	---	---

0	1	2	3	4	5	6	7	8
0	1	2	3	4	5	6	7	8
0	1	2	3	4	5	6	7	8

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro



pkt 2 timeout

send pkt2

record ack4 arrived

record ack5 arrived

Q: what happens when ack2 arrives?

0	1	2	3	4	5	6	7	8
0	1	2	3	4	5	6	7	8



Connection-oriented Transport TCP

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

› **point-to-point:**

- one sender, one receiver

› **reliable, in-order byte stream**

› **pipelined:**

- TCP congestion set window size

› **full duplex data:**

- bi-directional data flow in same connection

Assignment Project Exam Help - MSS: maximum segment

<https://eduassistpro.github.io/>

ion-oriented:

Add WeChat edu_assist_pro

ing (exchange of sgs) inits sender, receiver state before data exchange

› **flow controlled:**

- sender will not overwhelm receiver

TCP segment structure

URG: urgent data
(generally not used)

ACK: ACK #
valid

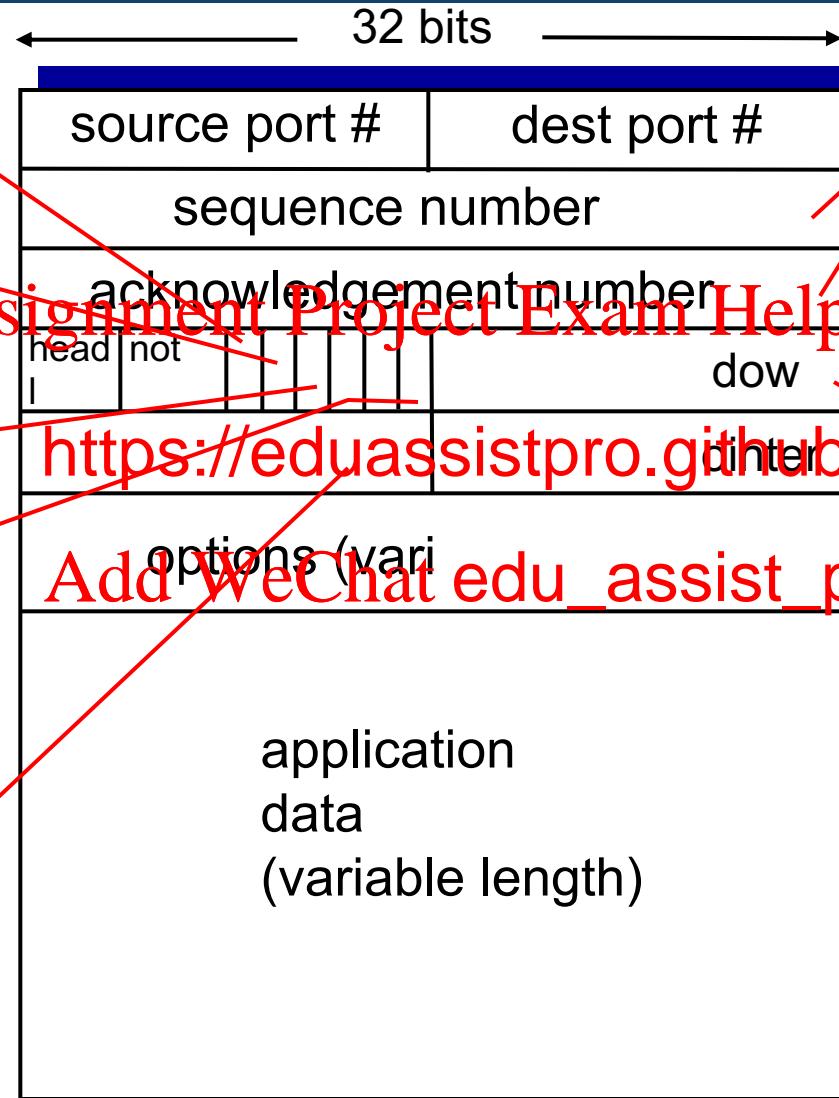
PSH: push data now
(generally not used)

RST, SYN, FIN:
connection estab
(setup, teardown
commands)

Internet
checksum
(as in UDP)

counting
by bytes
of data
(not segments!)

bytes
rcvr willing
to accept



sequence numbers:

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

acknowledgements Assignment Project Exam Help

- seq # of **next** by expected from ot <https://eduassistpro.github.io/>

- cumulative ACK

Q: how receiver handles out-of-order segments

- A: TCP spec doesn’t say,
- up to implementor
- Most will store, but still use cumulative ACK

outgoing segment from sender

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

Window size
 N



sent
ACKed

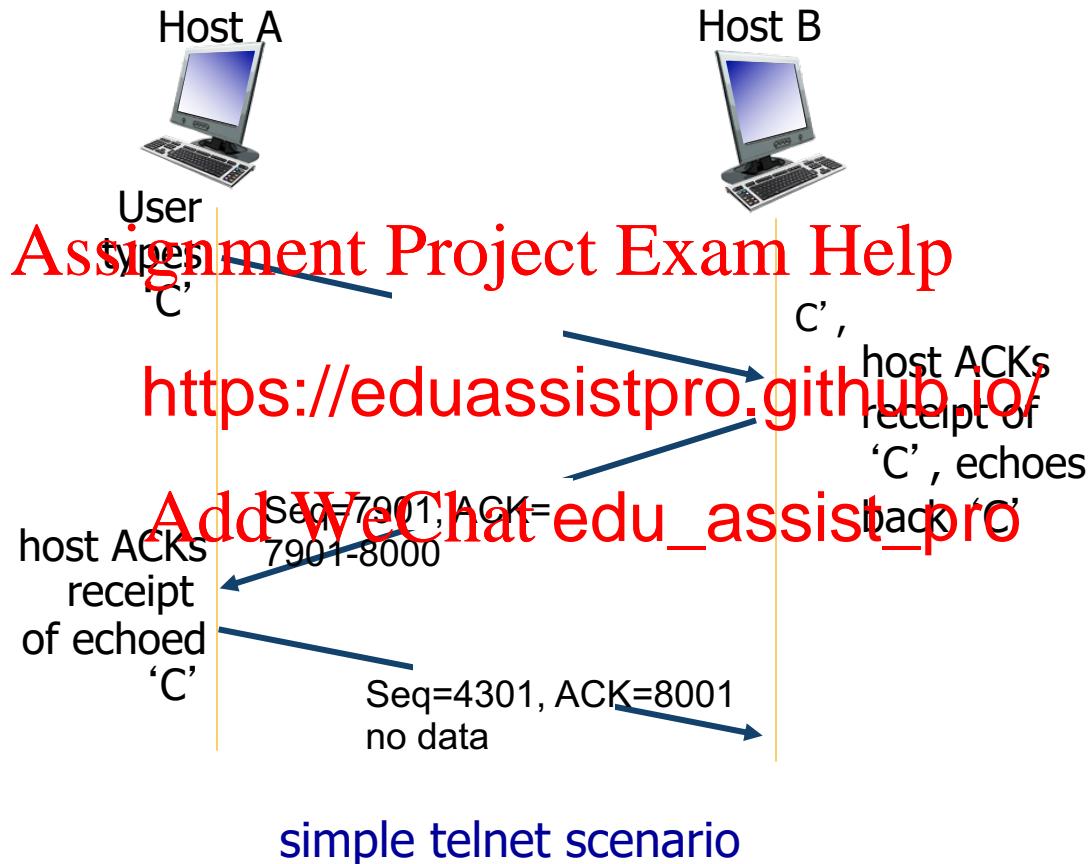
sent, not-
yet ACKed
("in-
flight")

usable
but not
yet sent

not
usable

incoming segment to sender

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

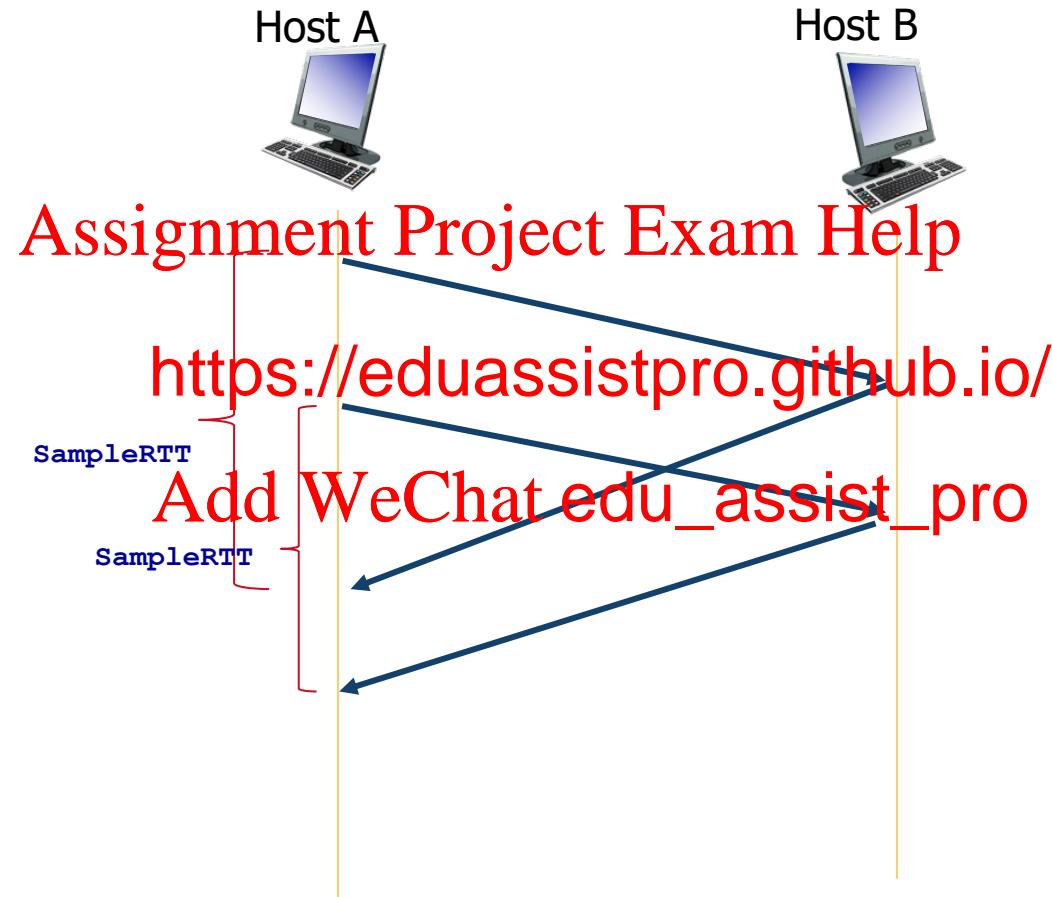


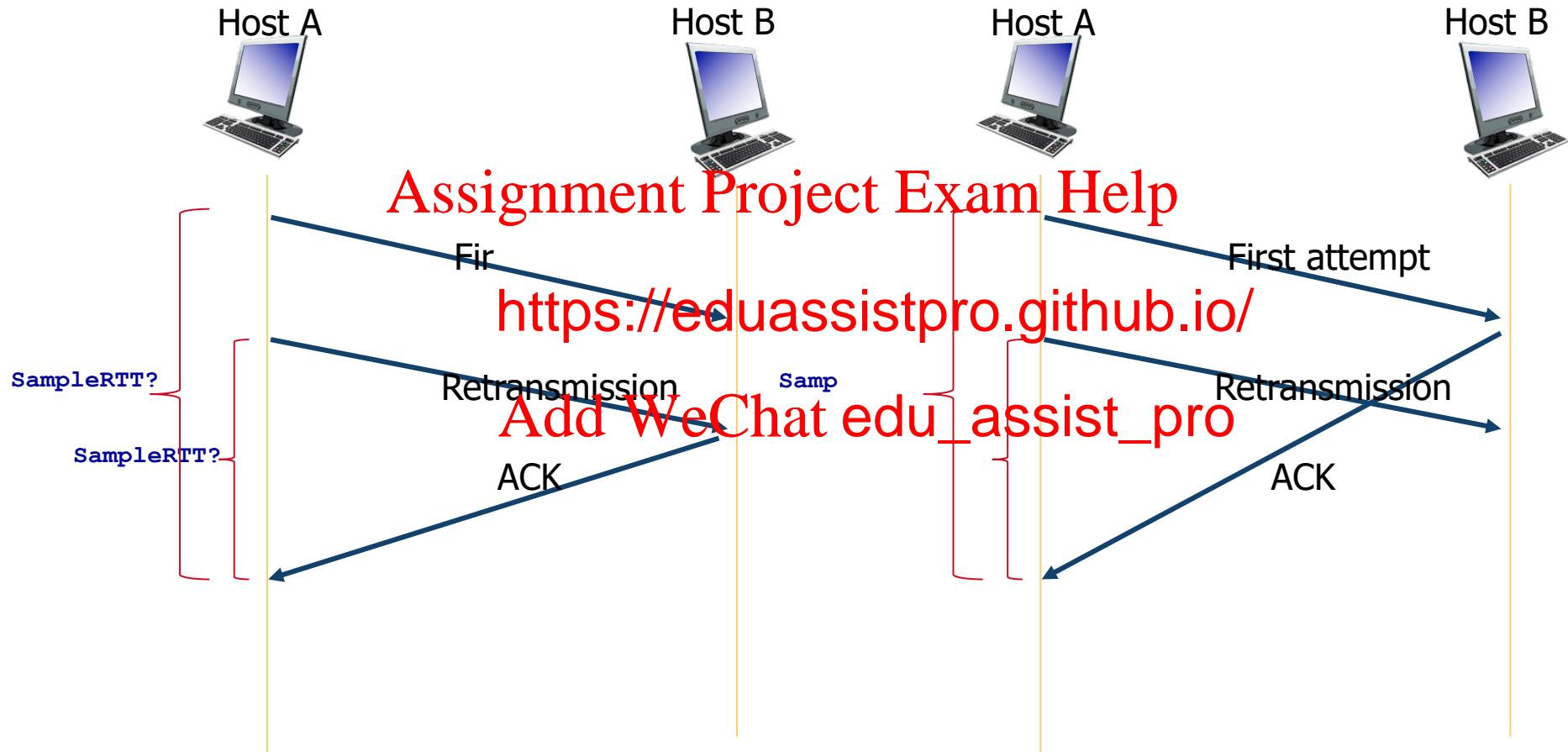
Q: how to set TCP timeout value?

- › 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 receipt
- › <https://eduassistpro.github.io/>
- › Add WeChat `edu_assist_pro`, will vary, want estimated RTT “smoother”
 - weighted average of several recent measurements, not just current **SampleRTT**

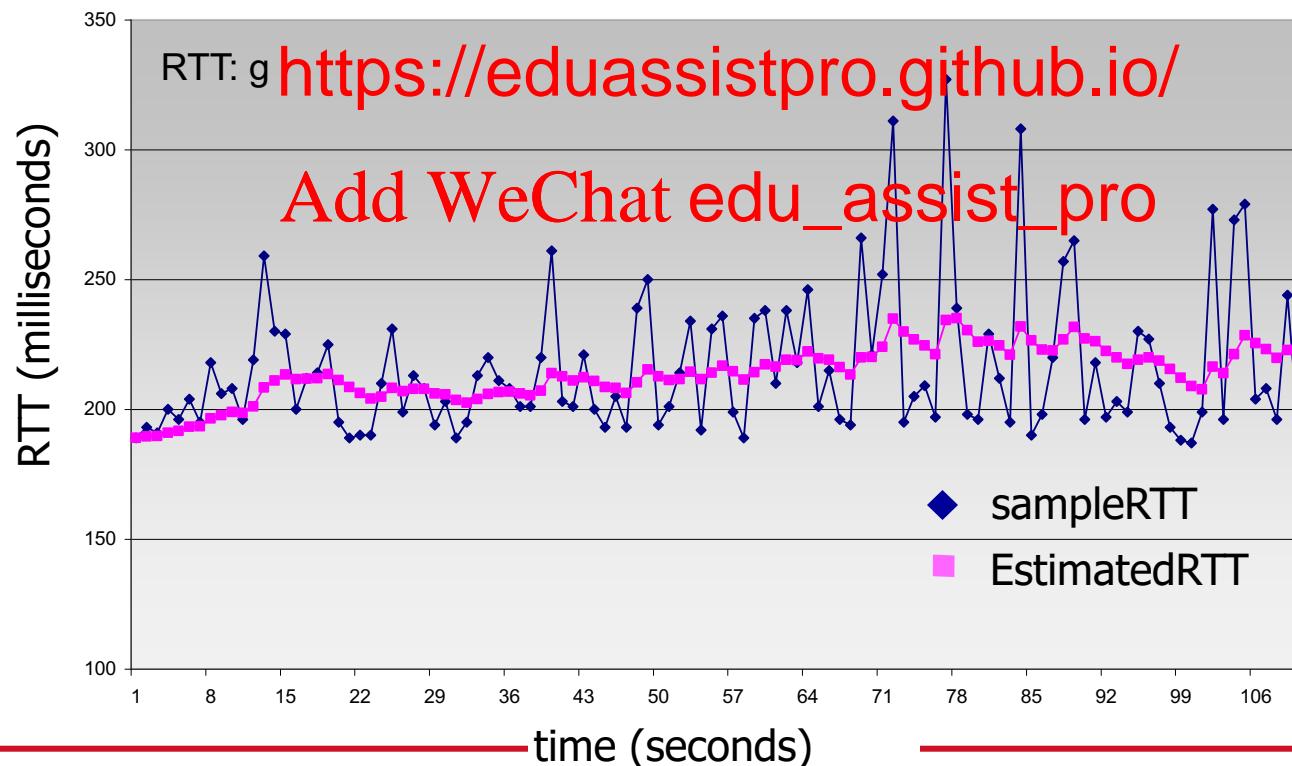






$$\text{EstimatedRTT} = (1 - \alpha) * \text{EstimatedRTT} + \alpha * \text{SampleRTT}$$

- ❖ exponential weighted moving average
- ❖ influence of past sample decreases exponentially fast
- ❖ typical value $\alpha = 0.25$



- › **timeout interval:** **EstimatedRTT** plus “safety margin”
 - large variation in **EstimatedRTT** -> larger safety margin
- › estimate SampleRTT deviation from EstimatedRTT:

DevRTT = (

<https://eduassistpro.github.io/>

(typically
Add WeChat edu_assist_pro)

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



↑
estimated RTT

↑
“safety margin”



Reliable Data Transfer in TCP

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro



- › TCP creates rdt service on top of IP's unreliable service

Assignment Project Exam Help

- pipelined segments
 - cumulative acks
 - single retransmission timer
- ly consider
TCP sender:
Add WeChat edu_assist_pro
- › retransmissions triggered by:
 - timeout events
 - duplicate acks
 - ignore flow control, congestion control

data rcvd from app:

- › create segment with seq #

› seq # is byte-stream number of first data byte in segment

- › start timer if no running

- think of timer as for oldest unacked segment
- expiration interval:
`TimeOutInterval`

timeout:

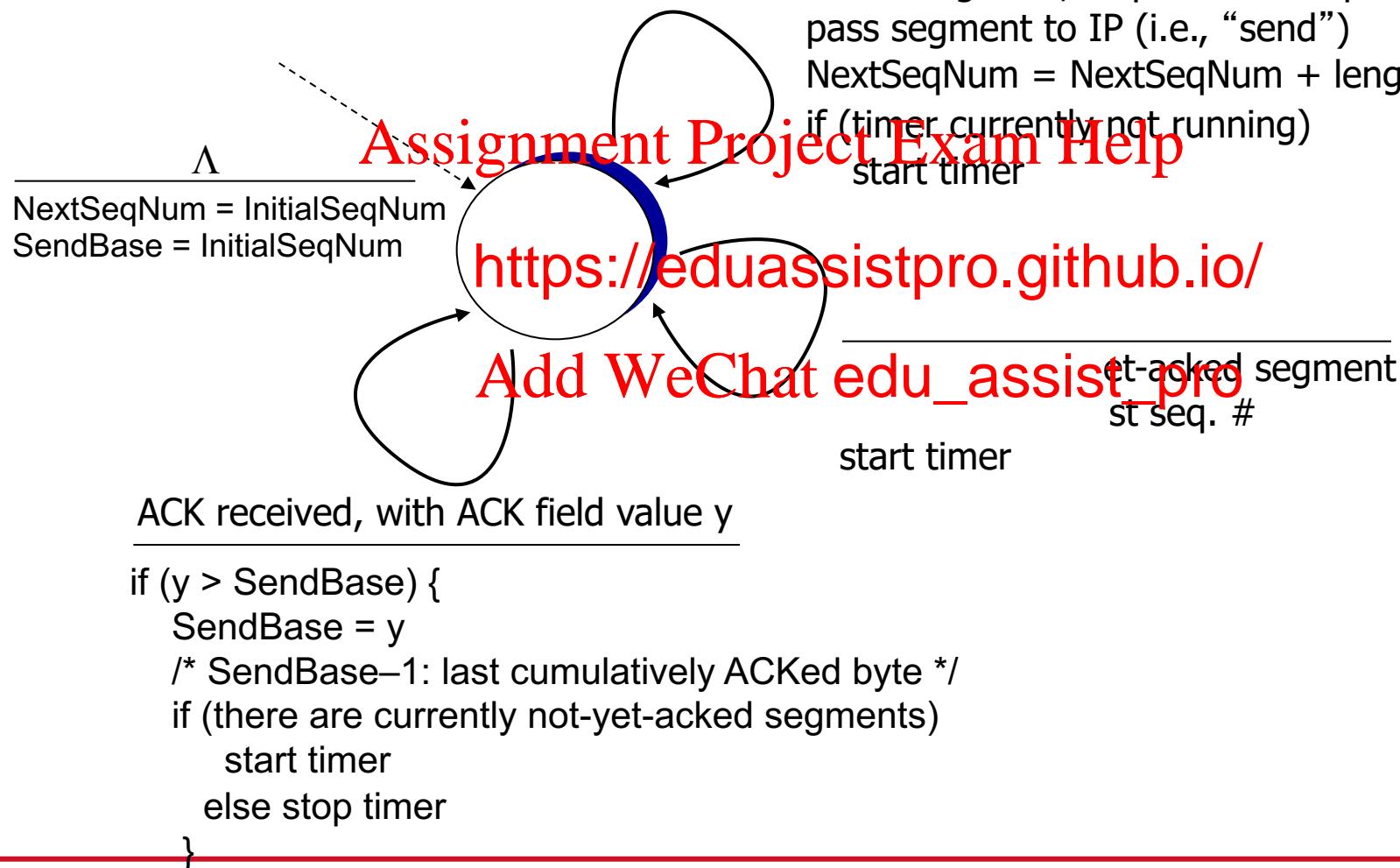
- › retransmit segment that caused timeout

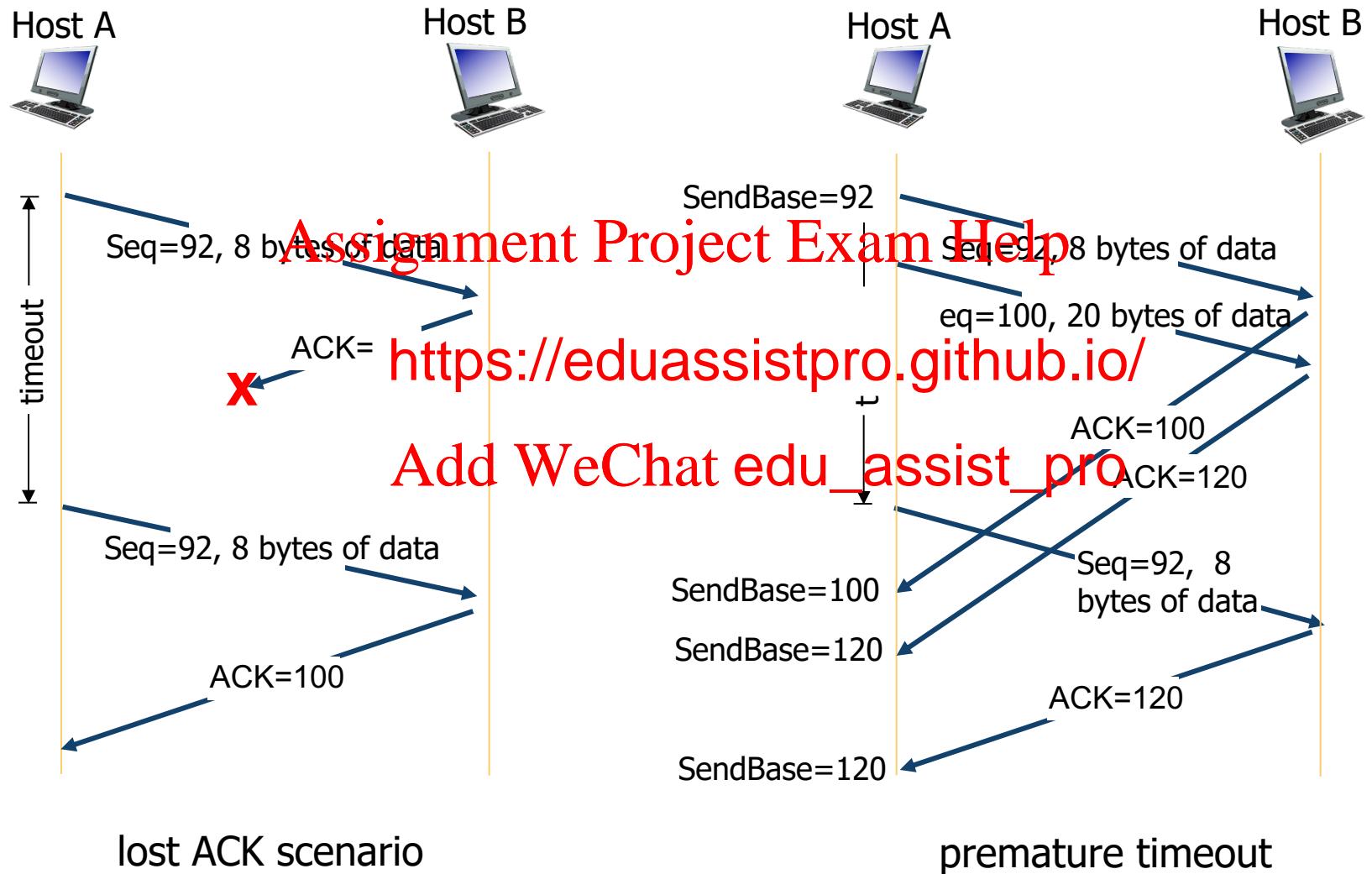
› restart timer

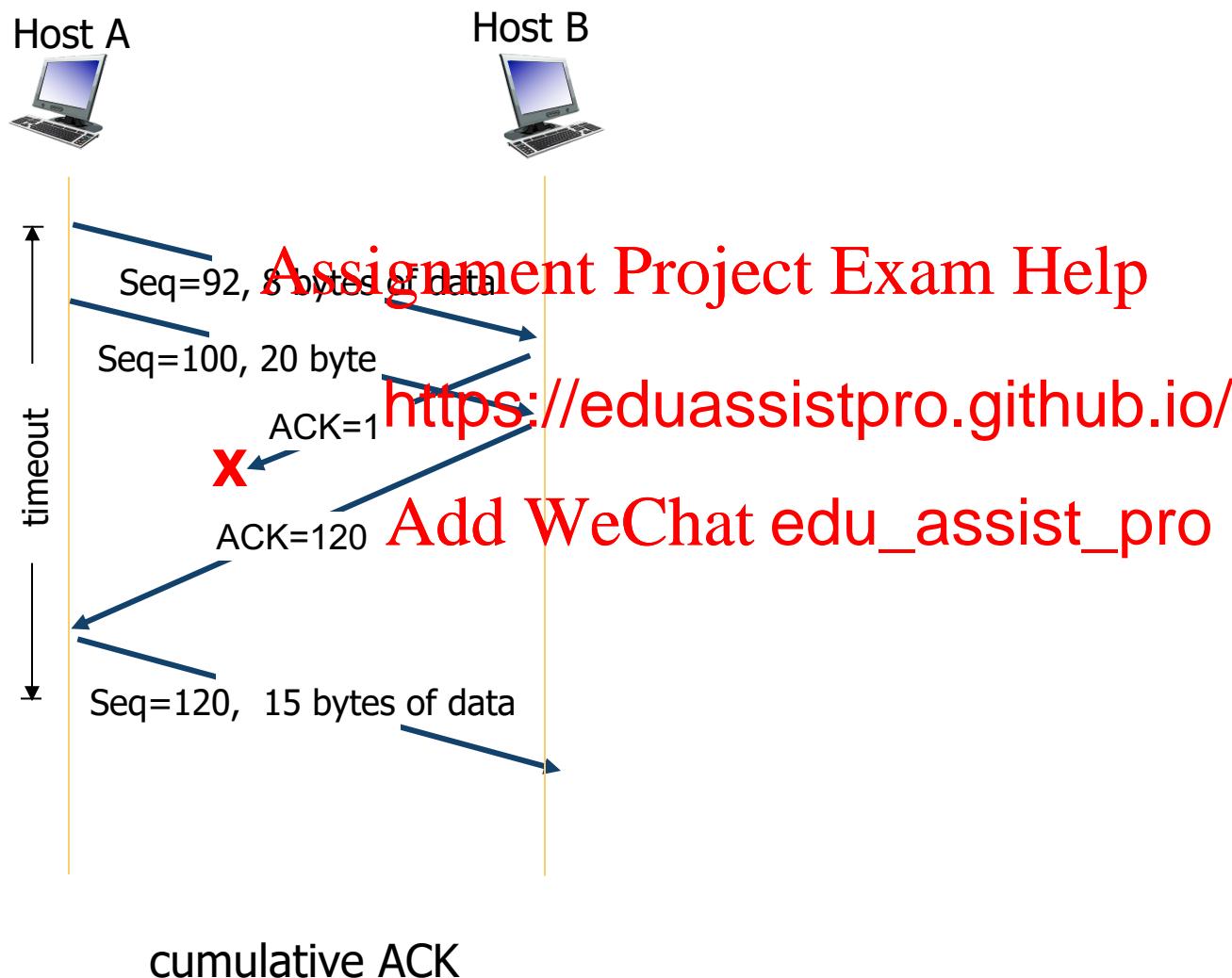
<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`

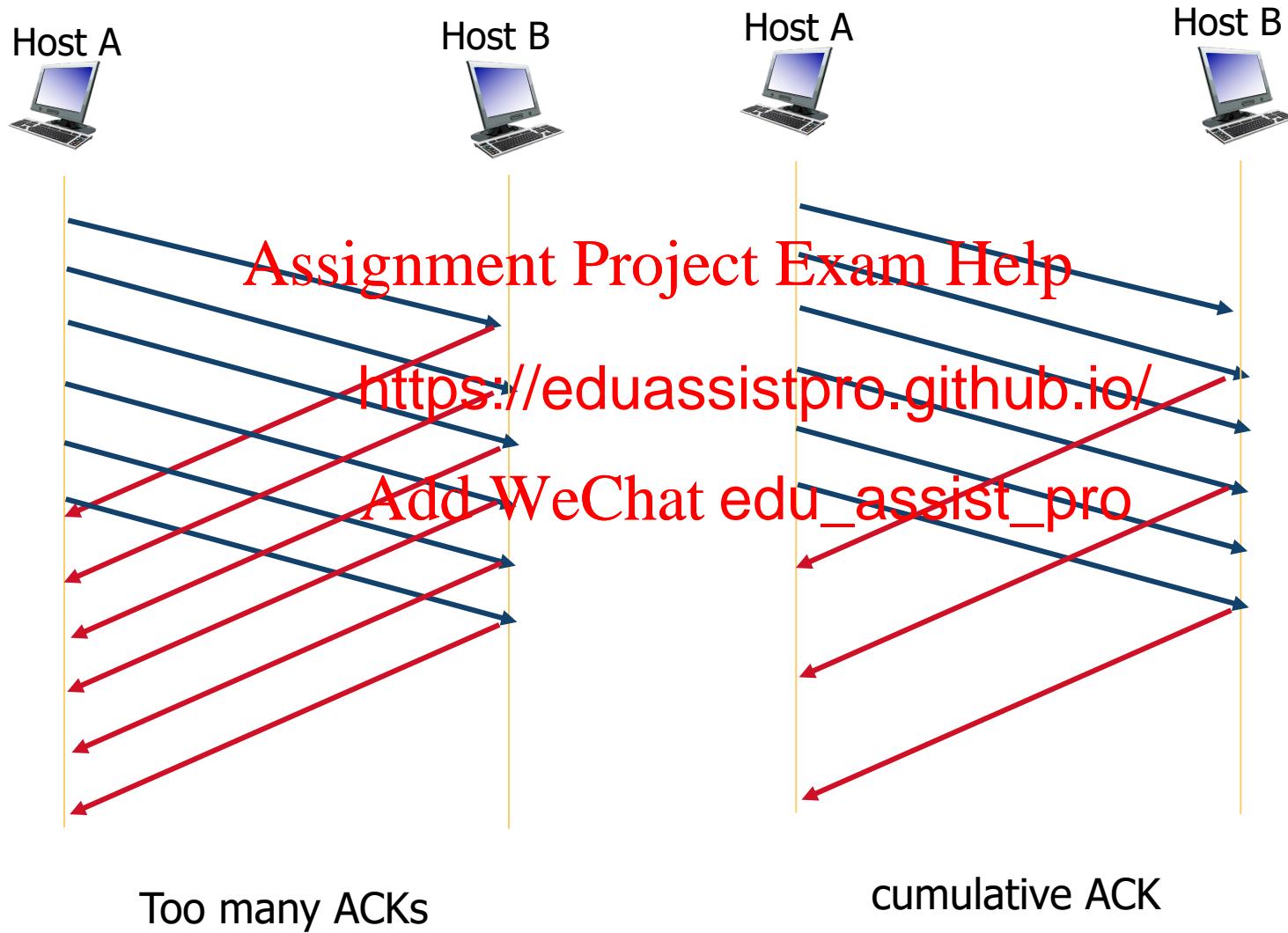
- update what is known to be ACKed
- start timer if there are still unacked segments

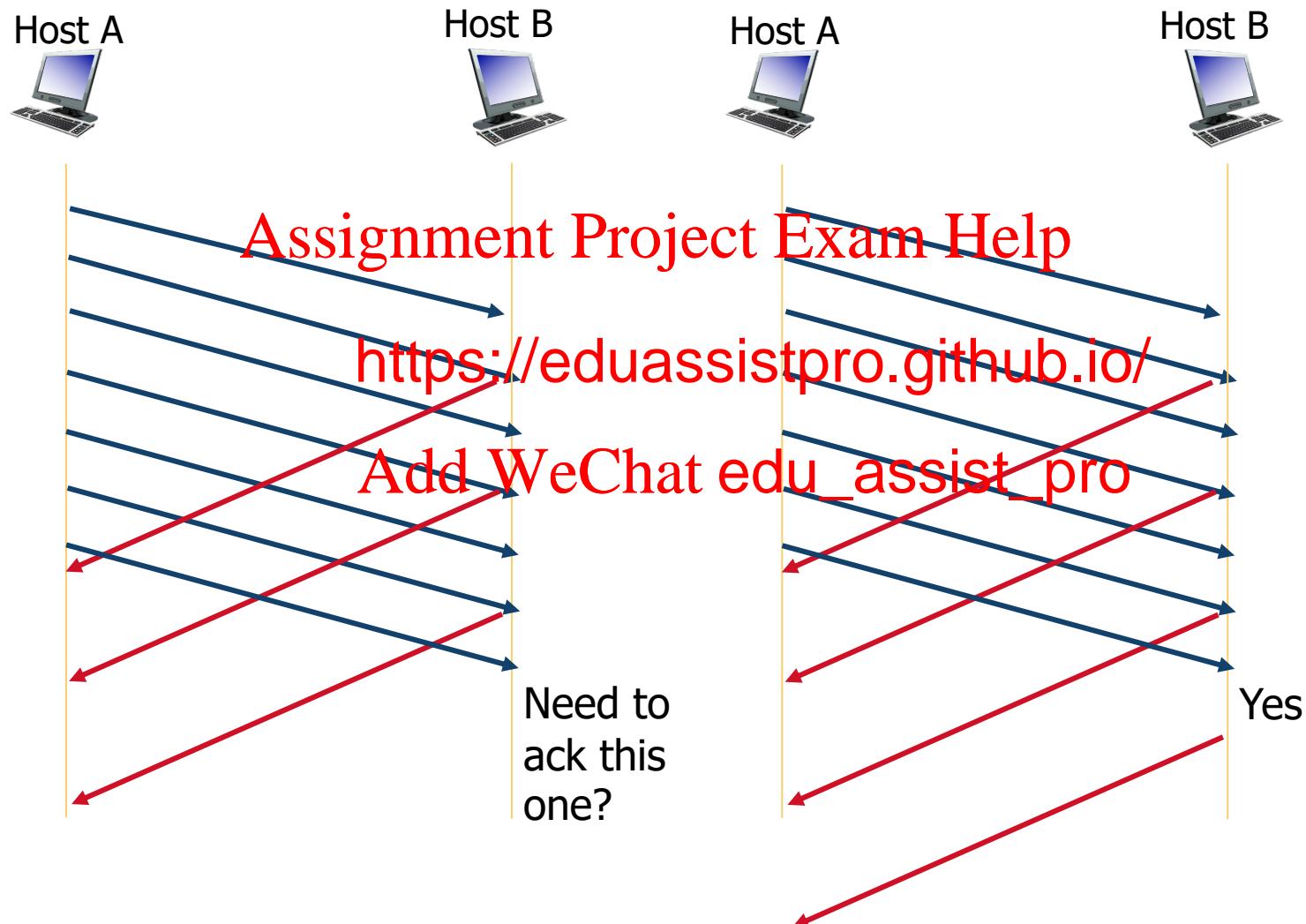






<i>event at receiver</i>	<i>TCP receiver action</i>
arrival of in-order segment with expected seq #. All data up to expected seq # already ACKed	delayed ACK. Wait up to 500ms for next segment. If no next segment, send ACK
arrival of in-order segment with expected seq #. One other segment has ACK pending	Send single cumulative ACK for all in-order segments up to the n^{th} in-order segments
arrival of out-of-order segment higher-than-expect seq. # . Gap detected	immediately send <i>duplicate ACK</i> , 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





- › time-out period often relatively long:

- long delay before resending lost packet

- › detect lost seg
- detect lost seg

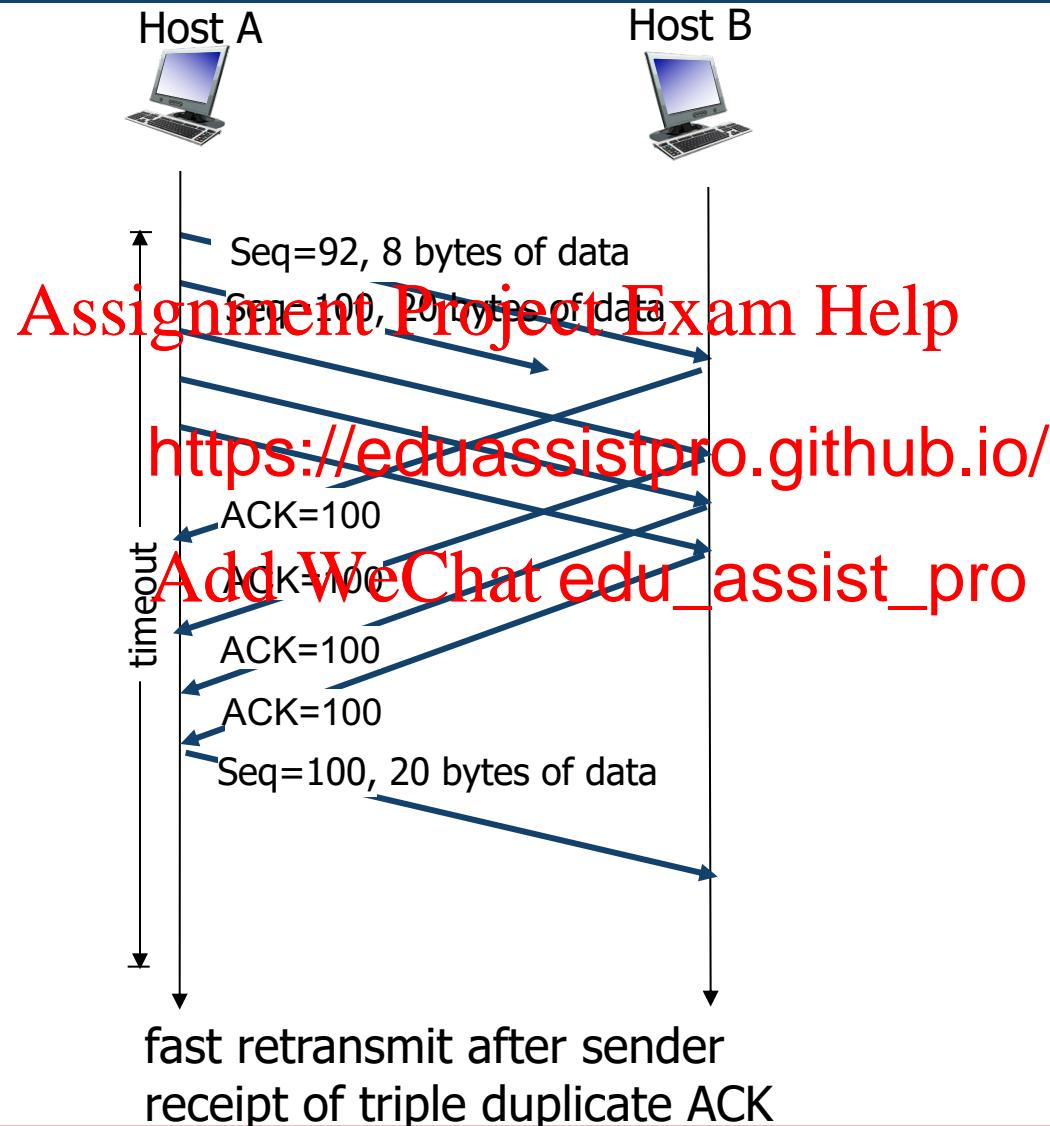
<https://eduassistpro.github.io/>

TCP fast retransmit

if sender receives 3
duplicate ACKs for same

- sender often sends many segments back-to-back
 - if segment is lost, there will likely be many duplicate ACKs.

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





Assignment Project Exam Help Flow Control in TCP

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

application may
remove data from
TCP socket buffers

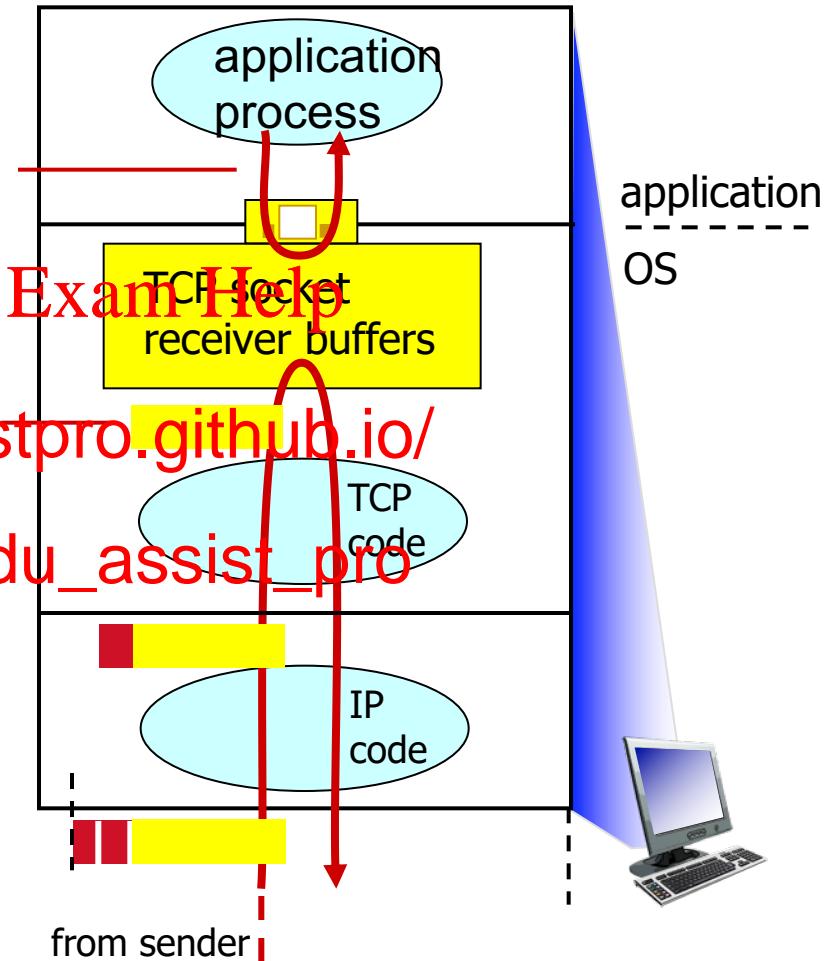
Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

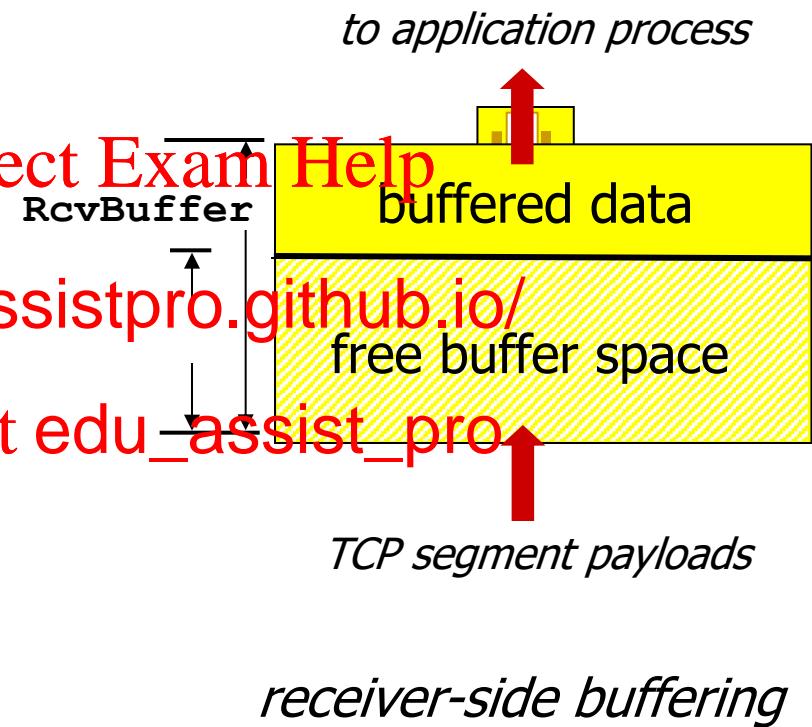
flow control

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



receiver protocol stack

- › receiver “advertises” free buffer space by including **rwnd** value in TCP header of receiver-to-sender segments
 - **RcvBuffer** size set options (typical defa <https://eduassistpro.github.io/> 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





Connection Management in TCP

Assignment Project Exam Help

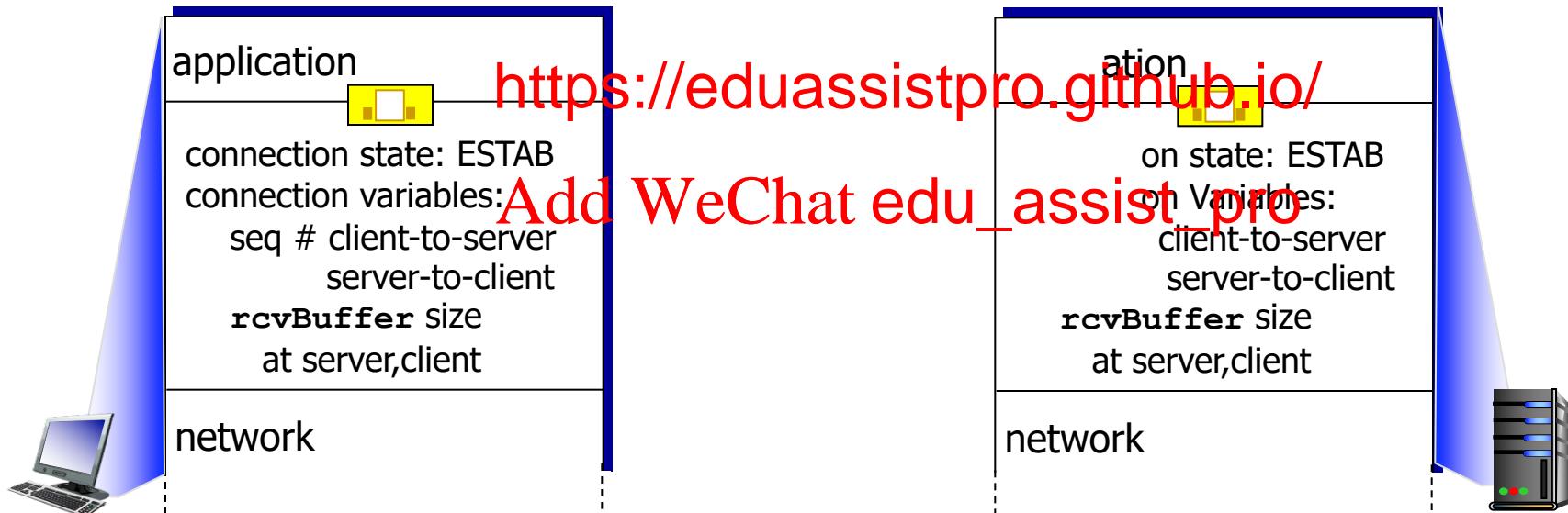
<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

before exchanging data, sender/receiver “handshake”:

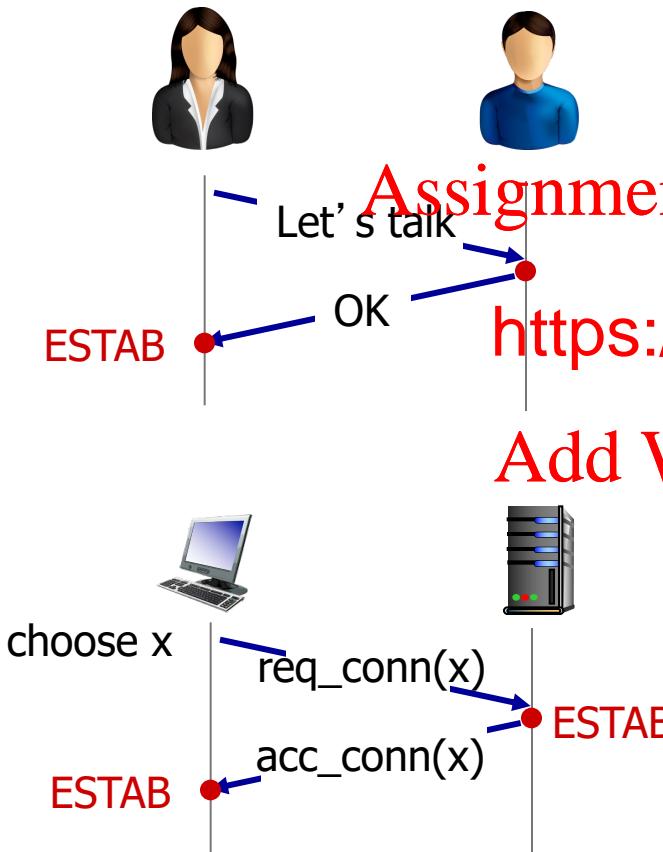
- › agree to establish connection (each knowing the other willing to establish connection)
- › agree on common parameters

Assignment Project Exam Help



Agreeing to establish a connection

2-way handshake:



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

variable delays
 retransmitted messages (e.g.
 due to message loss
 dering)

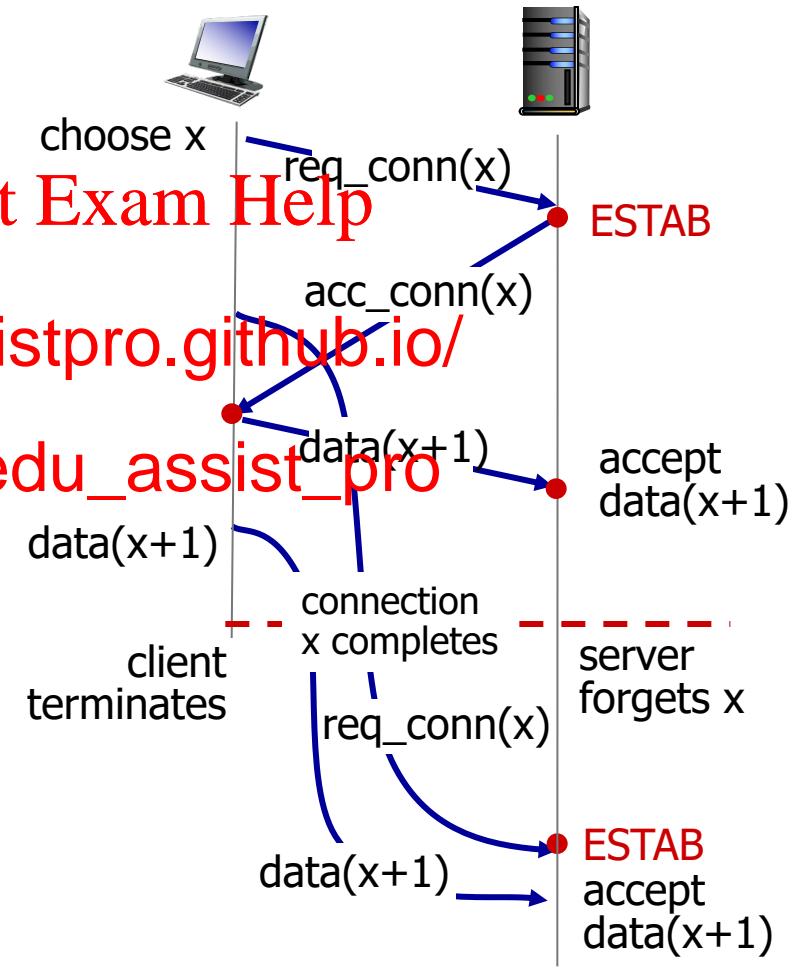
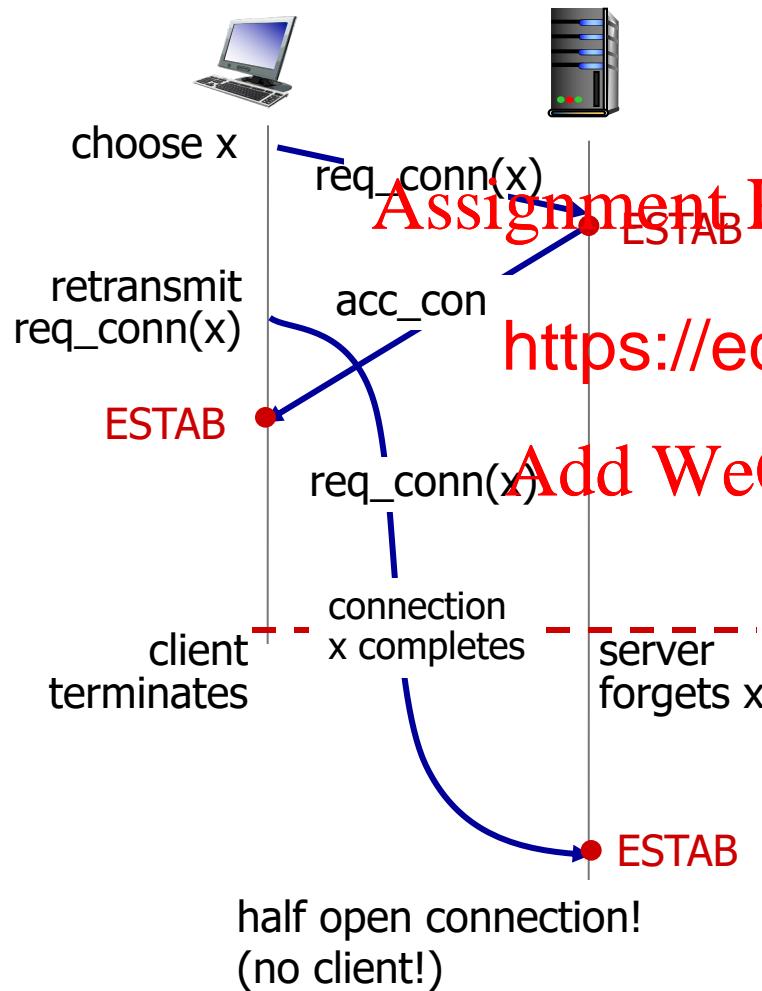
Assignment Project Exam Help

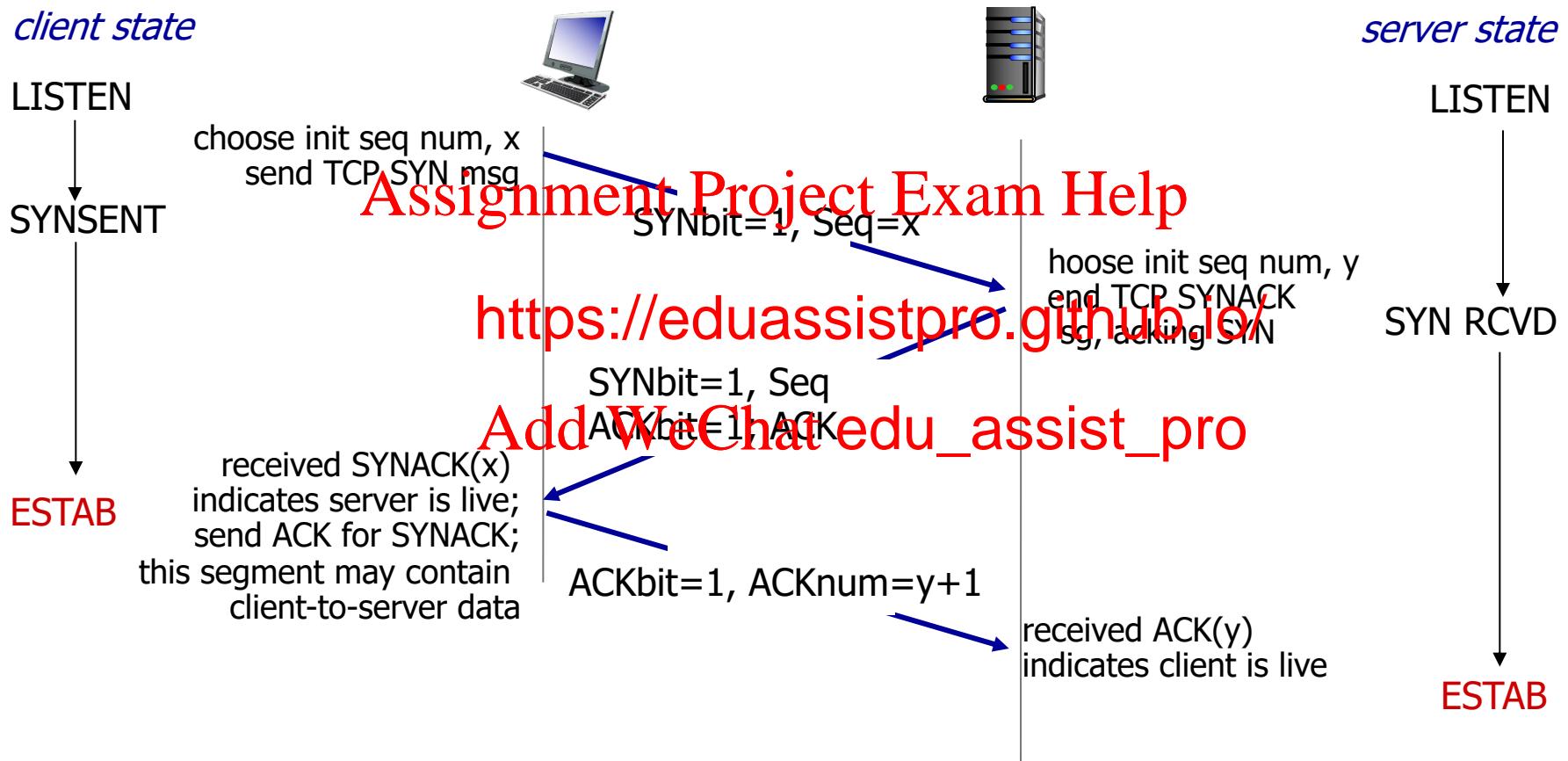
<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Agreeing to establish a connection

2-way handshake failure scenarios:







- › client, server each closes their side of connection
 - send TCP segment with FIN bit = 1
- › respond to received FIN with ACK

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

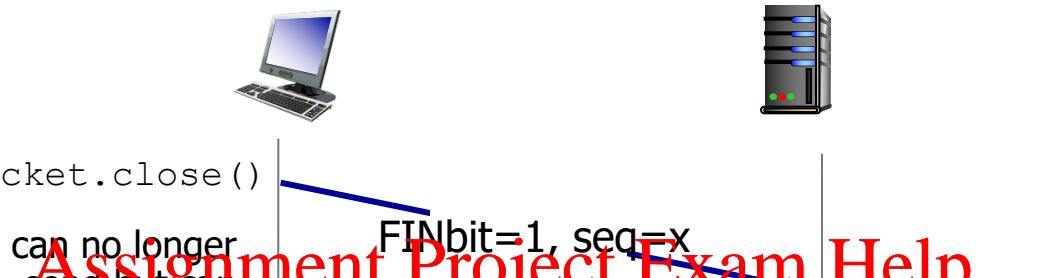
client state

ESTAB
↓
FIN_WAIT_1

FIN_WAIT_2

TIMED_WAIT

CLOSED



wait for <https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

FINbit=1,
ACKbit=1; ACKnum=y+1

timed wait for $2 \times \text{max segment lifetime}$

server state

ESTAB
↓
CLOSE_WAIT

LAST_ACK

CLOSED



can still send data

can no longer send data

TCP segment structure

URG: urgent data
(generally not used)

ACK: ACK #
valid

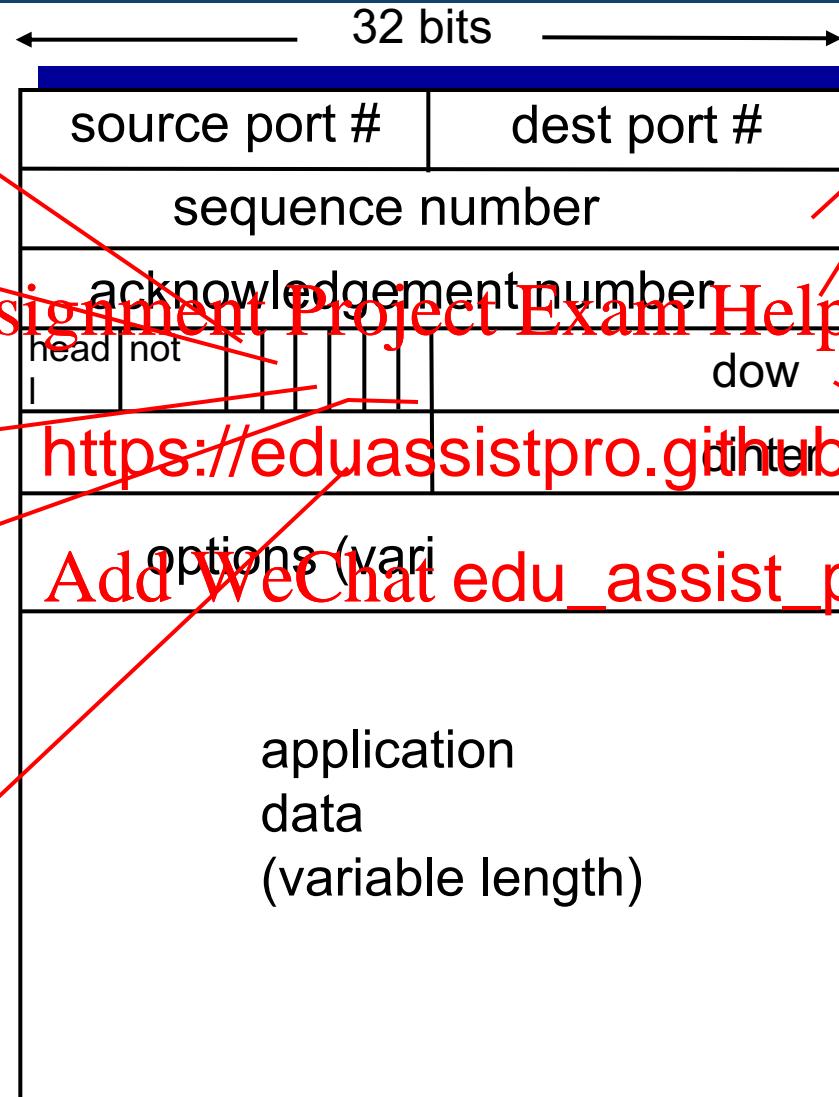
PSH: push data now
(generally not used)

RST, SYN, FIN:
connection estab
(setup, teardown
commands)

Internet
checksum
(as in UDP)

counting
by bytes
of data
(not segments!)

bytes
rcvr willing
to accept





Principles of Congestion Control

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

congestion:

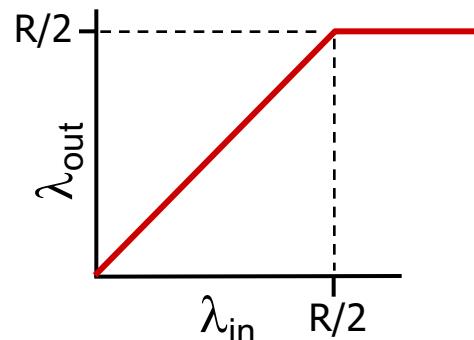
- › informally: “too many sources sending too much data too fast fo **Assignment Project Exam Help**
- › different from f <https://eduassistpro.github.io/>
- › manifestations: **Add WeChat edu_assist_pro**
 - lost packets (buffer overflow at routers)
 - long delays (queueing in router buffers)
- › a top-10 problem!

Causes/costs of congestion: scenario 1

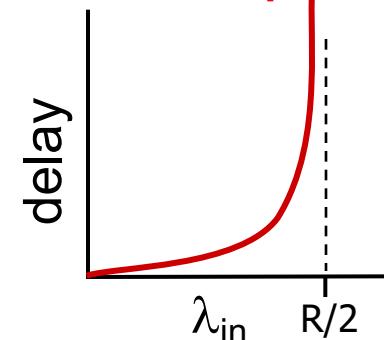
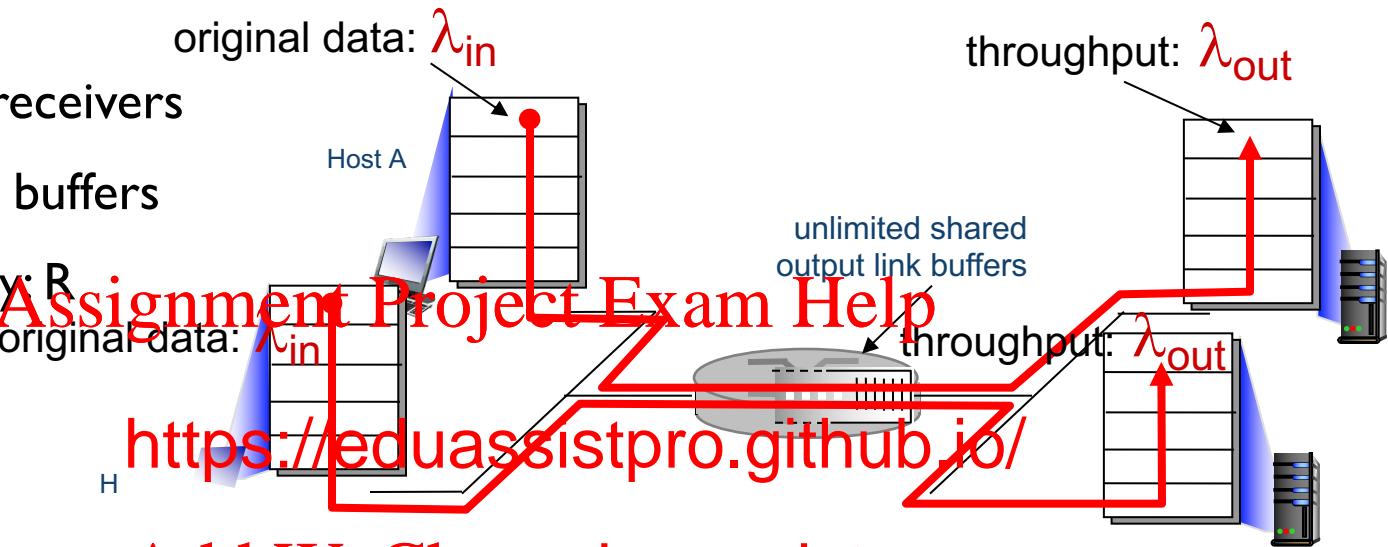
- › two senders, two receivers
- › one router, infinite buffers
- › output link capacity: R
- › no retransmission

Assignment Project Exam Help
<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`



- › maximum per-connection throughput: $R/2$

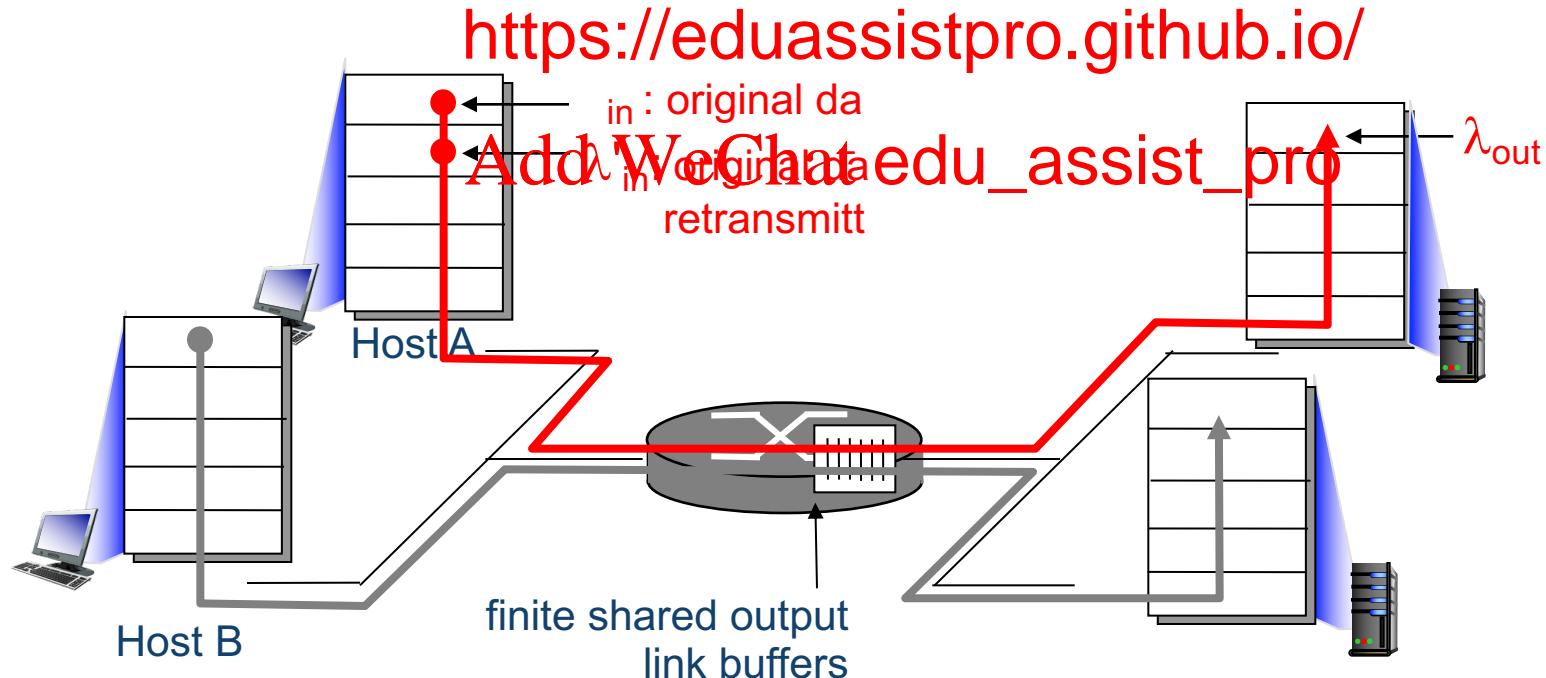


- ❖ large delays as arrival rate, λ_{in} , approaches capacity

Causes/costs of congestion: scenario 2

- › one router, *finite buffers*
- › sender retransmission of timed-out packet
 - application-layer input = application-layer output: $\lambda_{in} = \lambda_{out}$,
 - Goodput
 - transport-layer input includes retransmissions: $\lambda_{in}^{TCP} \neq \lambda_{in}$

Assignment Project Exam Help

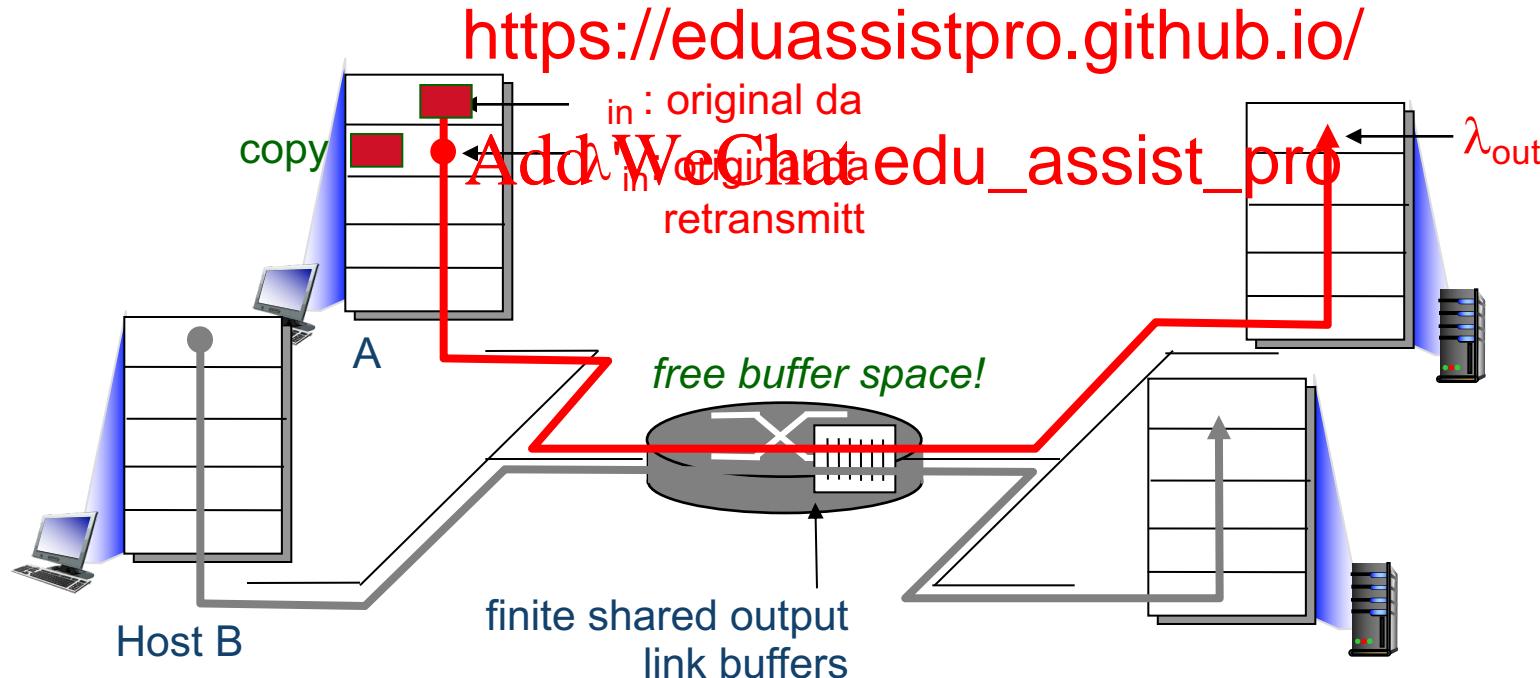
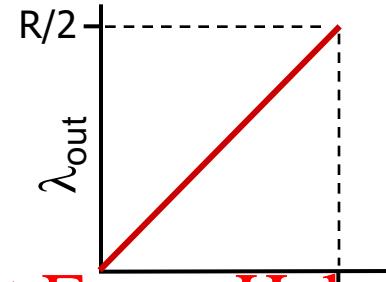


Causes/costs of congestion: scenario 2

idealization: perfect knowledge

- › sender sends only when router buffers available

Assignment Project Exam Help



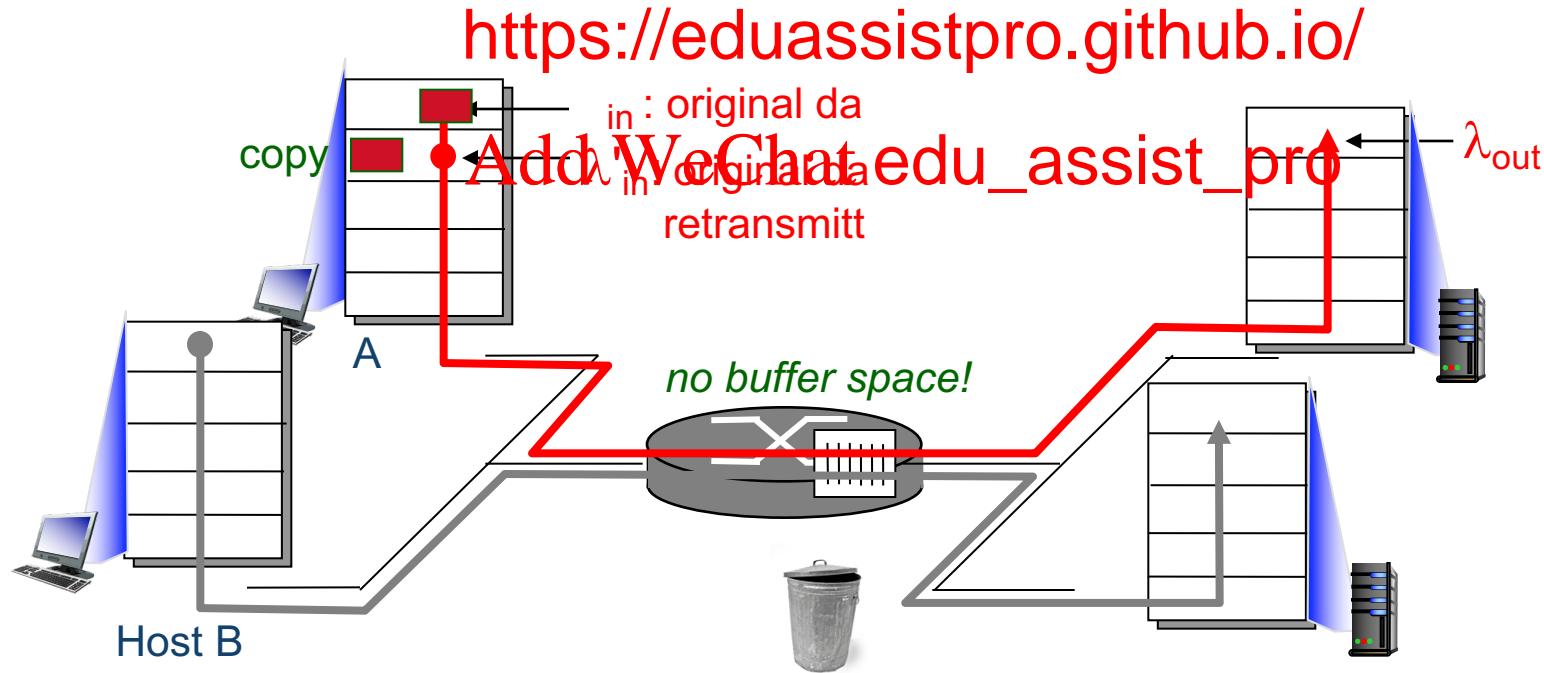
Causes/costs of congestion: scenario 2

Idealization: known loss

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

- › sender only resends if
packet known to be lost

[Assignment](#) [Project](#) [Exam](#) [Help](#)

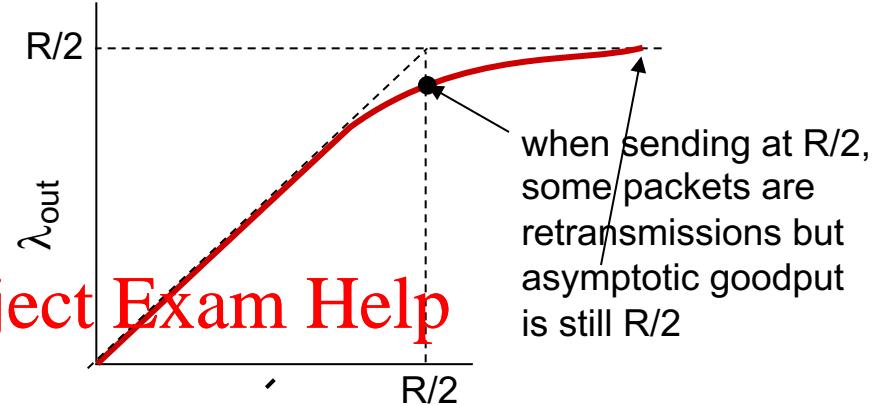


Causes/costs of congestion: scenario 2

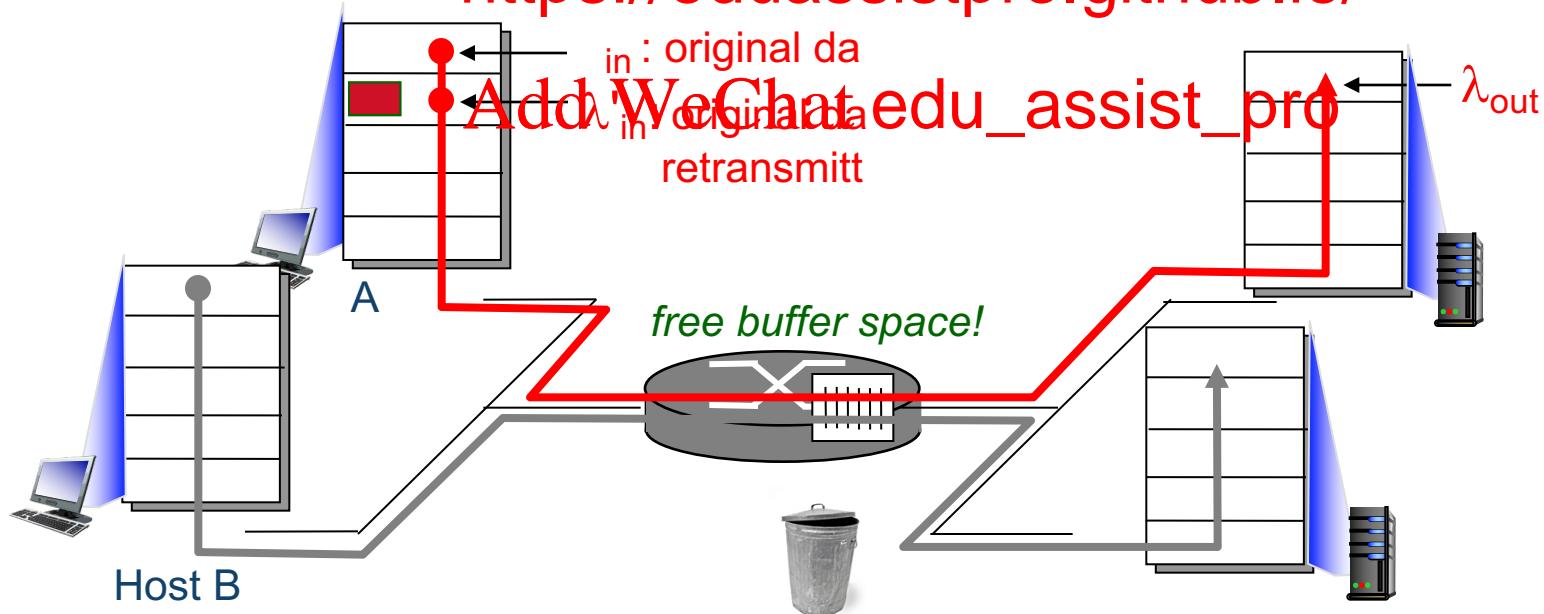
Idealization: known loss

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

- › sender only resends if packet known to be lost



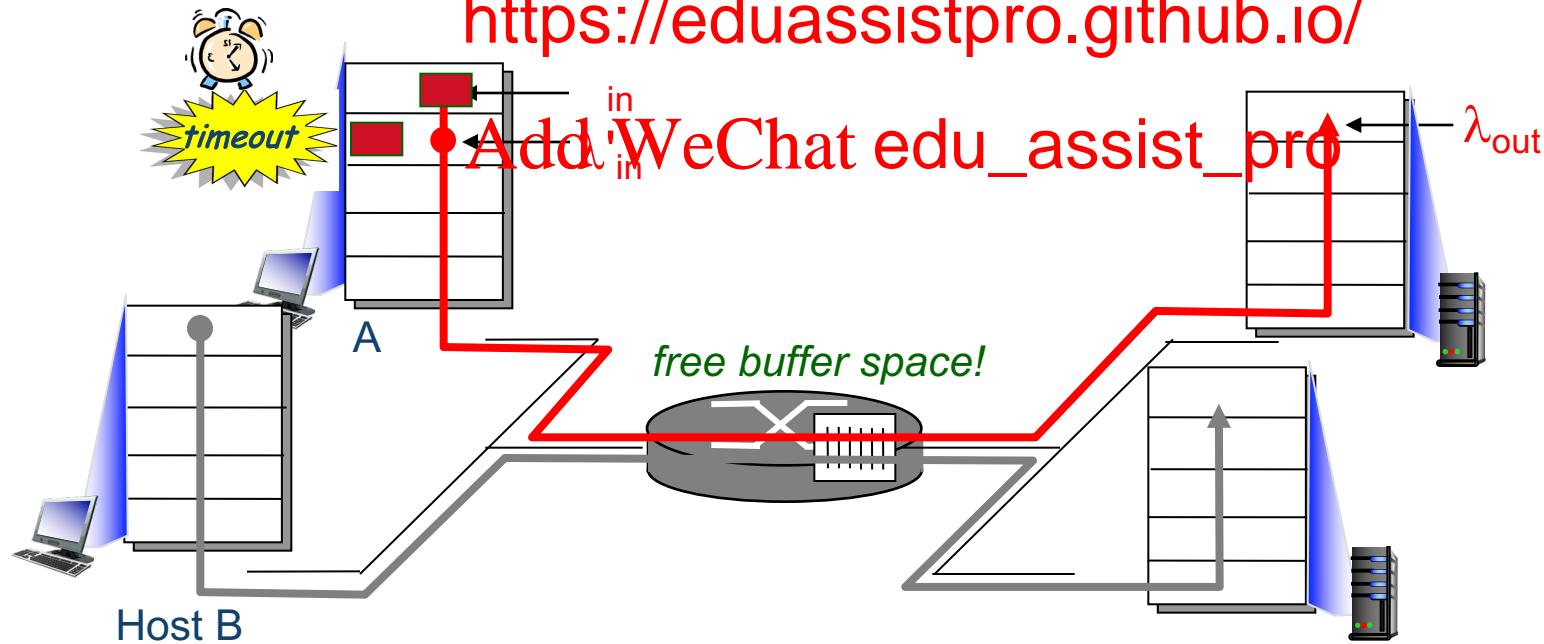
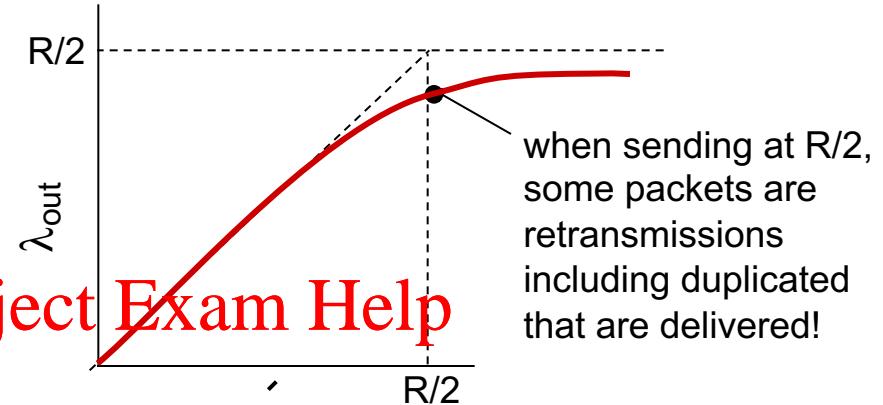
<https://eduassistpro.github.io/>



Causes/costs of congestion: scenario 2

Realistic: *duplicates*

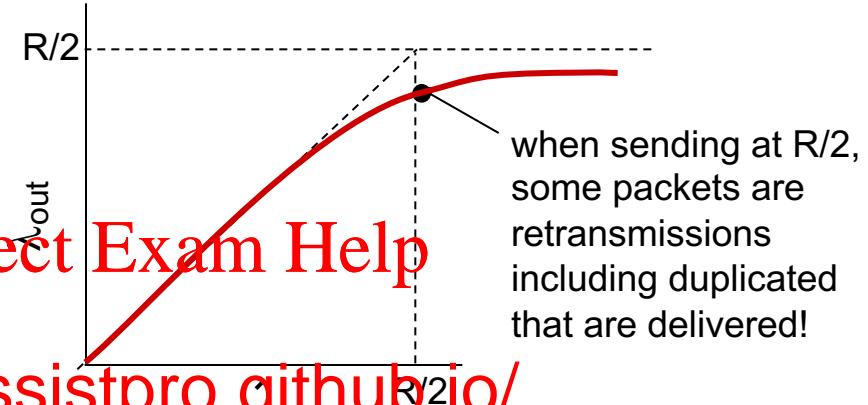
- ❖ packets can be lost, dropped at router due to full buffers
- ❖ sender times out prematurely, sending **two** copies, both of which are delivered



Realistic: *duplicates*

- ❖ packets can be lost, dropped at router due to full buffers
- ❖ sender times out prematurely, sending *two* copies which are delivered

Assignment Project Exam Help
<https://eduassistpro.github.io/>



“costs” of congestion: Add WeChat edu_assist_pro

- ❖ more work (retrans) for given “goodput”
- ❖ unneeded retransmissions: link carries multiple copies of pkt
 - decreasing goodput

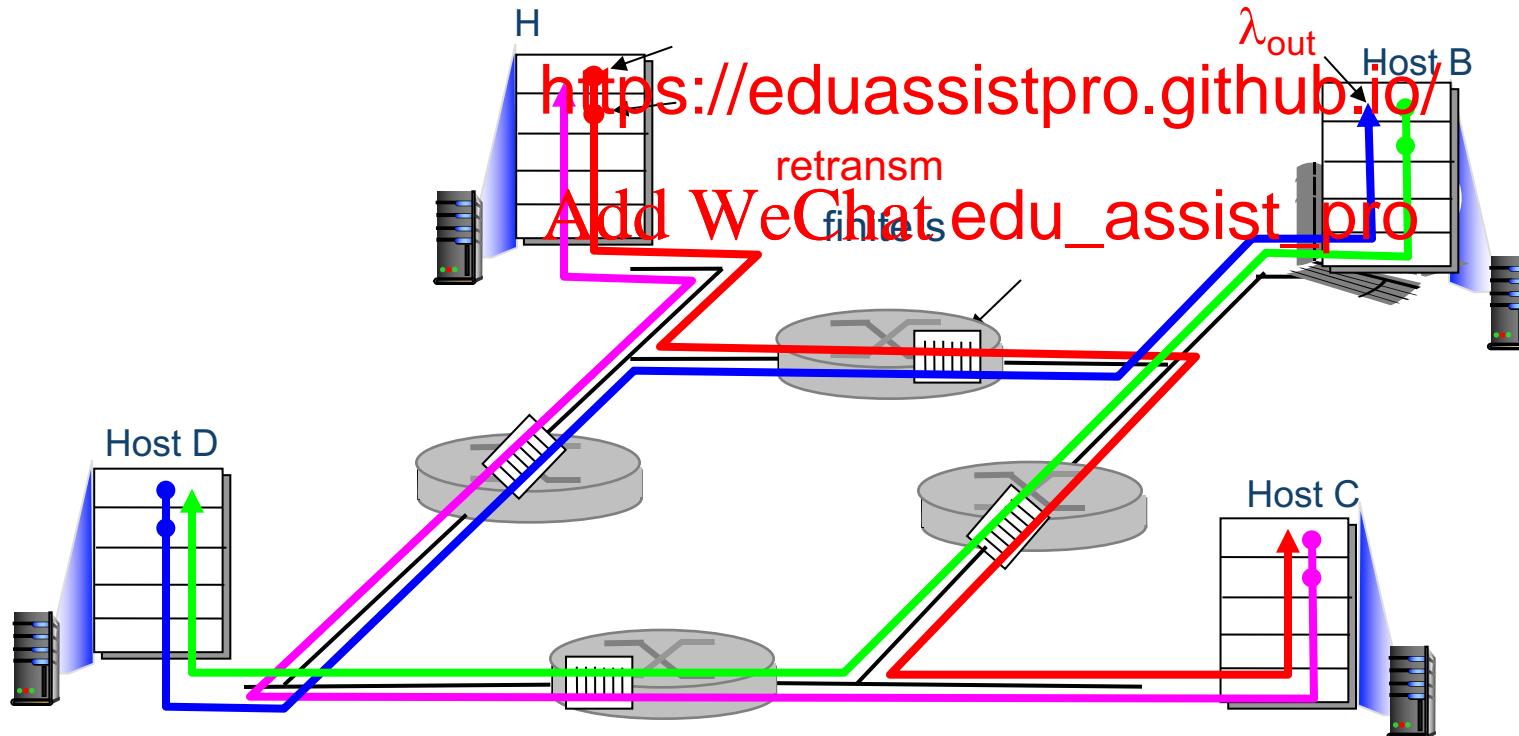
Causes/costs of congestion: scenario 3

- › four senders
- › multihop paths
- › timeout/retransmit

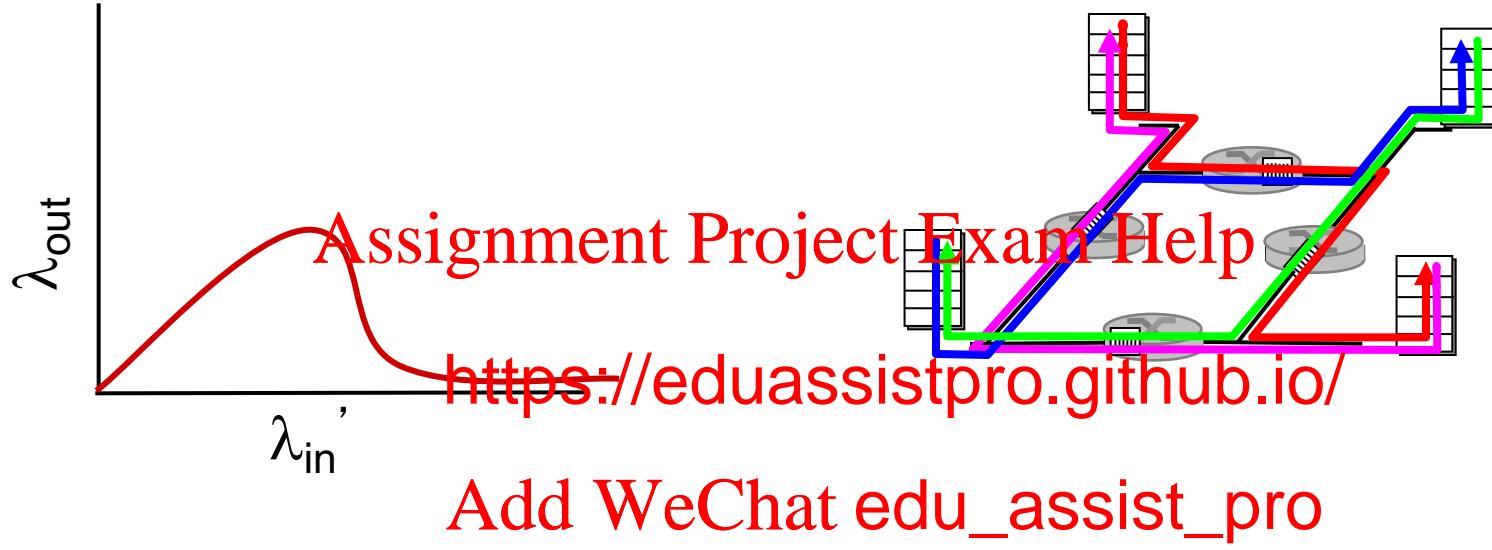
Q: what happens as λ_{in} ' increases ?

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

Assignment Project Exam Help



Causes/costs of congestion: scenario 3



another “cost” of congestion:

- ❖ when packet dropped, any “upstream transmission capacity used for that packet was wasted!

two broad approaches towards congestion control:

~~end-end congestion control~~ Project Firewall-Hopping

control:

- › no explicit feed from network
- › congestion inferred from end-system observed loss, delay
- › approach taken by TCP

<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`

tion control:

- single bit indicating congestion
- explicit rate for sender to send at



TCP Congestion Control

Assignment Project Exam Help

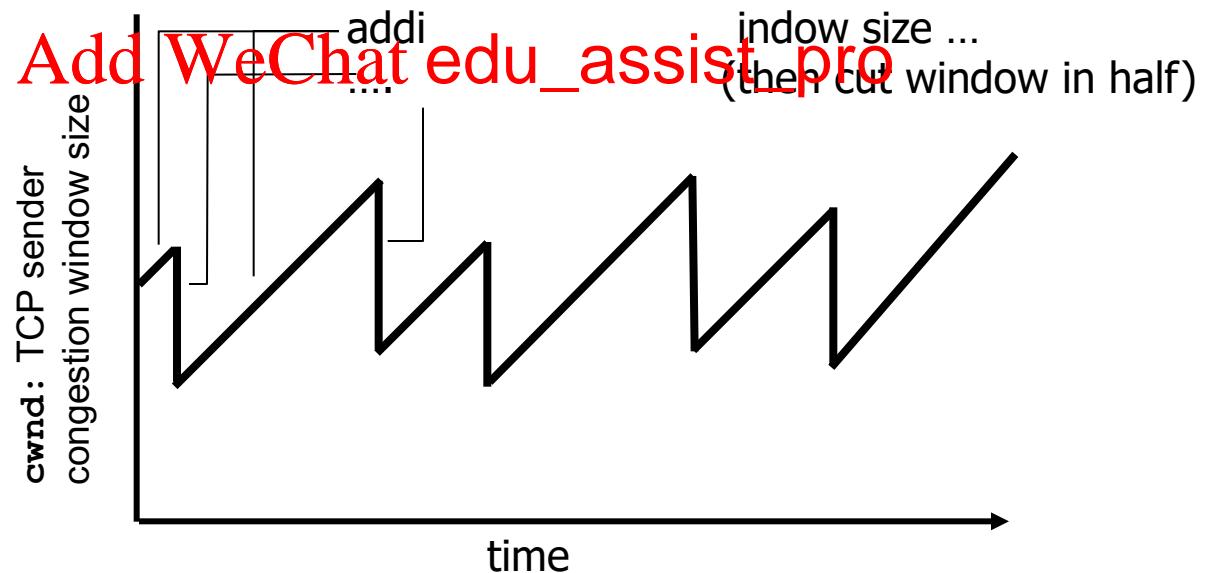
<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

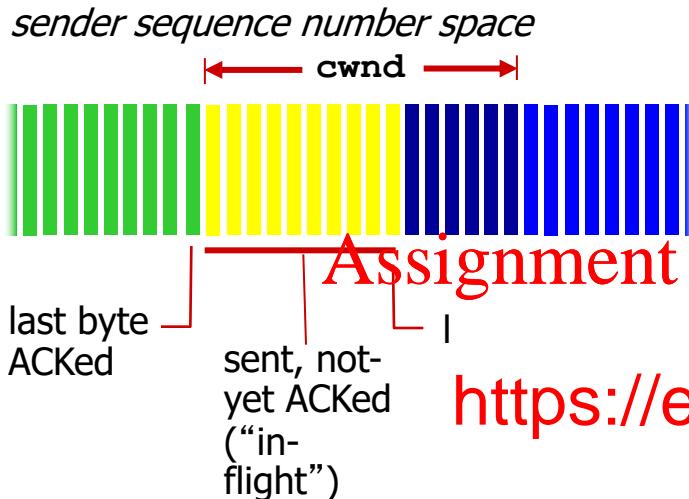
Additive increase multiplicative decrease (AIMD)

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

AIMD saw tooth behavior: probing for bandwidth



TCP Congestion Control: details



TCP sending rate:

- › roughly: send cwnd bytes, wait RTT for ACKS, then send more bytes

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro $\frac{\text{cwnd}}{\text{RTT}}$ bytes/sec

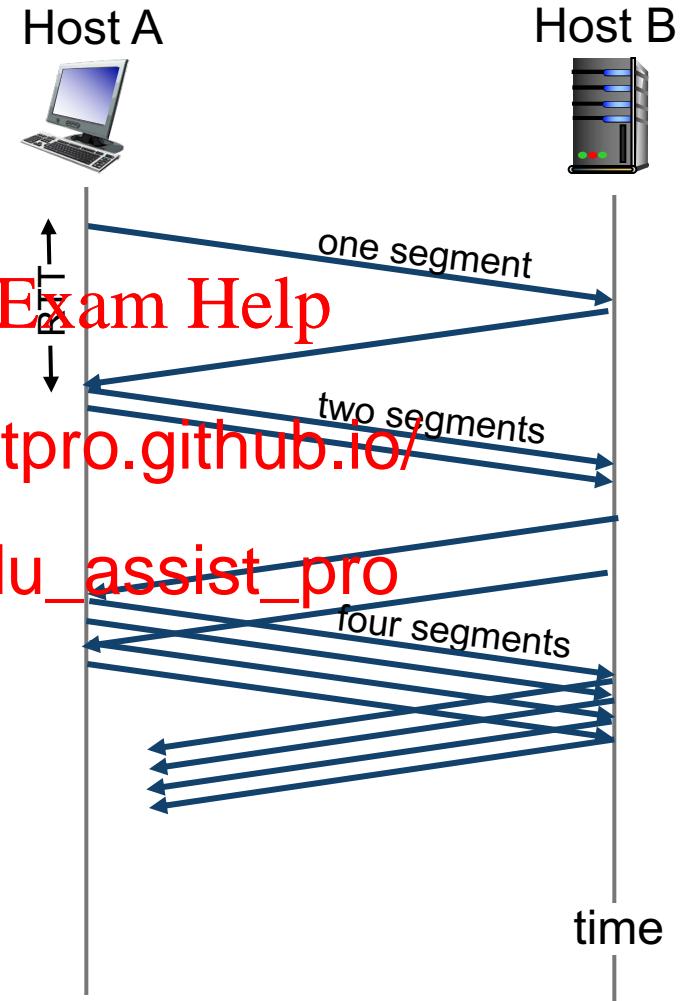
LastByteSent - LastByteAcked \leq cwnd

- › cwnd is dynamic, function of perceived network congestion

- › when connection begins, increase rate exponentially:

- initially $cwnd = 1$ MSS
- double $cwnd$ every RTT
- done by incrementing every ACK received

- › summary: initial rate is slow but ramps up exponentially fast
- › when should the exponential increase switch to linear (additive increase)?



TCP: switching from slow start to CA

Q: when should the exponential increase switch to linear?

A: cwnd reaches **ssthresh**

cwnd=12 ssthresh=6
loss!

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Implementation:

Add WeChat `edu_assist_pro`

- › At beginning **ssthresh**, specified in different versions of TCP
- › (In this example **ssthresh=8 segment**)
- › on loss event, **ssthresh** is set to 1/2 of **cwnd** just before loss event

› loss indicated by timeout:

- **cwnd** set to 1 MSS;
- window then grows exponentially (as in slow start) to **ssthresh**, then grows linearly

<https://eduassistpro.github.io/>

› loss indicated by

- › TCP Tahoe, same as loss indicated by 3 duplicate acks

› TCP RENO

- **cwnd** is cut in half window then grows linearly (additive increase)
- fast recovery



Assignment Project Exam Help

additive increase

additive increase

tiplicative decease

<https://eduassistpro.github.io/>

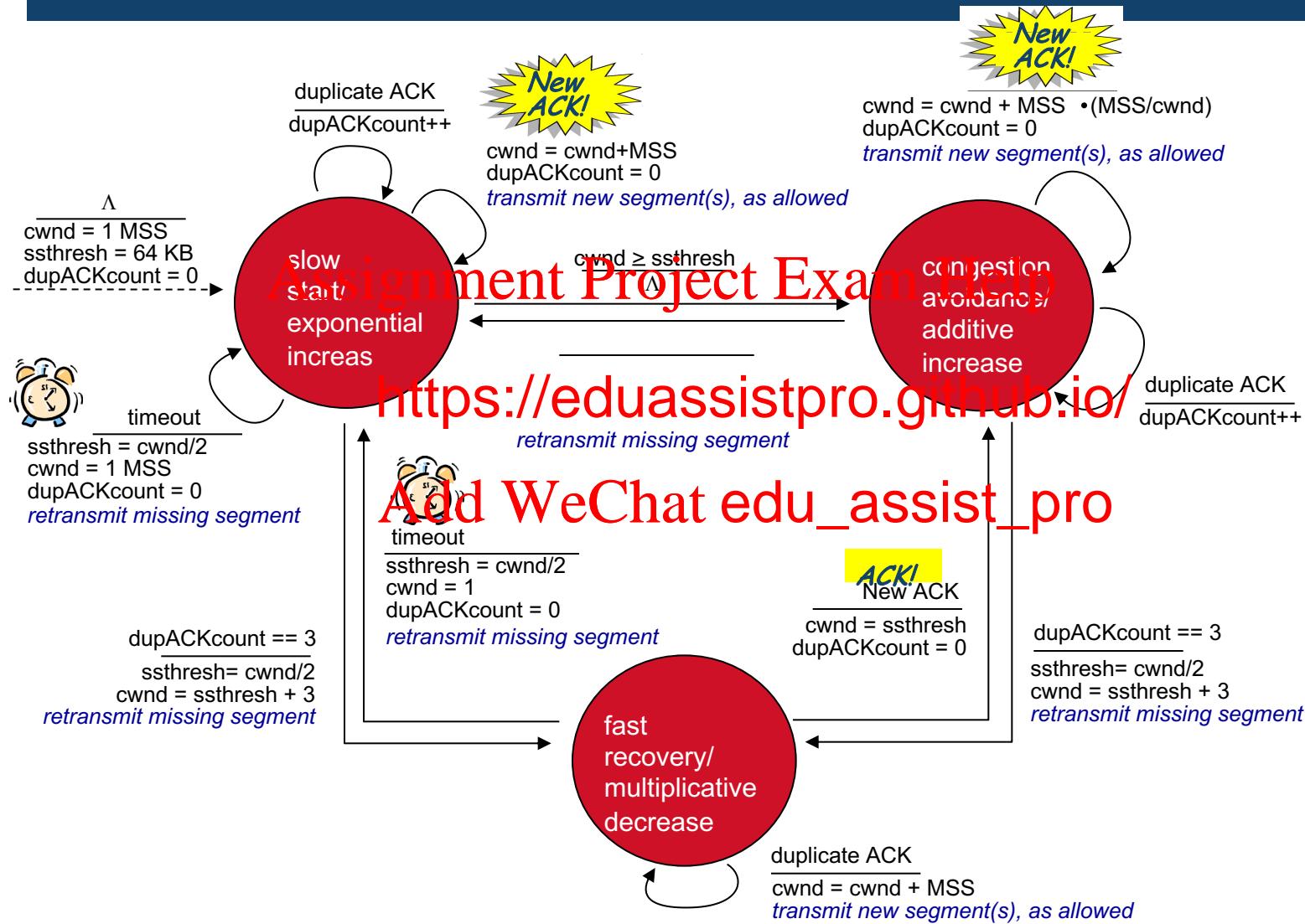
additive increase

slow start

Add WeChat edu_assist_pro

slow start

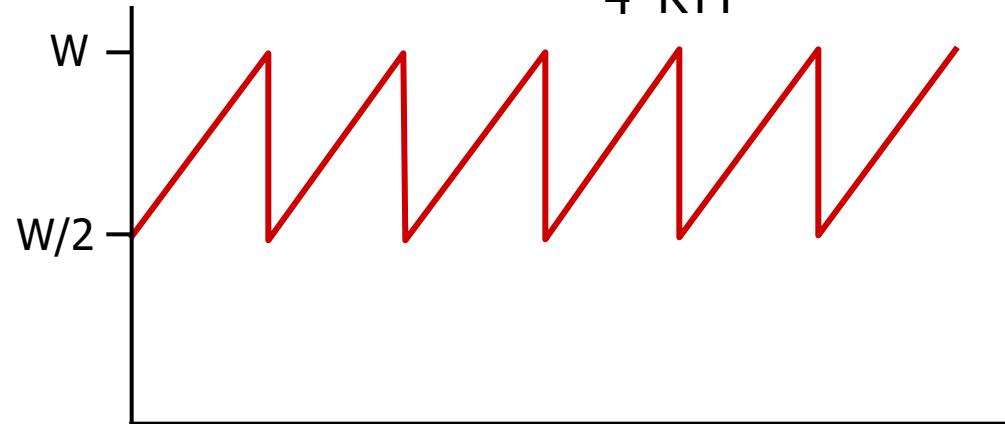
Summary: TCP Reno Congestion Control



- › avg. TCP thruput as function of window size, RTT?
 - ignore slow start, assume always data to send
- › W: windowAssignment Project Exam Help
- avg. window si <https://eduassistpro.github.io/>
- avg. thruput is $3/4W$ per RTT

Add WeChat edu_assist_pro

$$\text{avg TCP thruput} = \frac{3}{4} \frac{W}{\text{RTT}} \text{ bytes/sec}$$



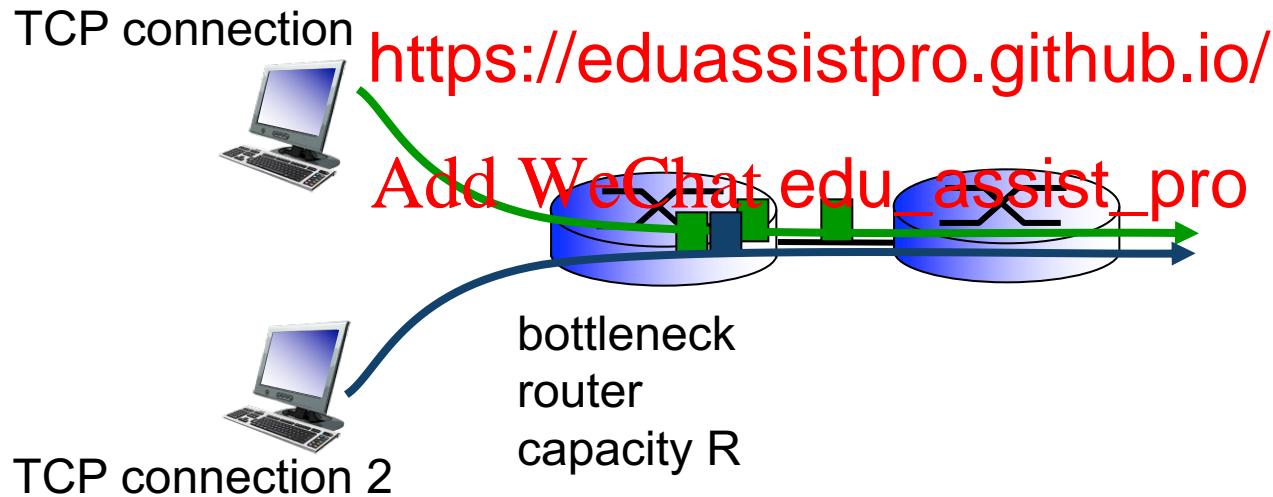
TCP Futures: TCP over “long, fat pipes”

- › example: 1500 byte segments, 100ms RTT, want 10 Gbps throughput
- › requires $W = 83,333$ in-flight segments
- › throughput in terms of $\frac{1}{L}$ [Mathis 1997]:
$$\text{TCP throughput} = \frac{1}{L}$$

<https://eduassistpro.github.io/>
- to achieve 10 Gbps throughput, need a loss rate of $L = 2 \cdot 10^{-10}$
– *a very small loss rate!*
- › new versions of TCP for high-speed
 - › Vegas, Westwood, CUBIC, etc.

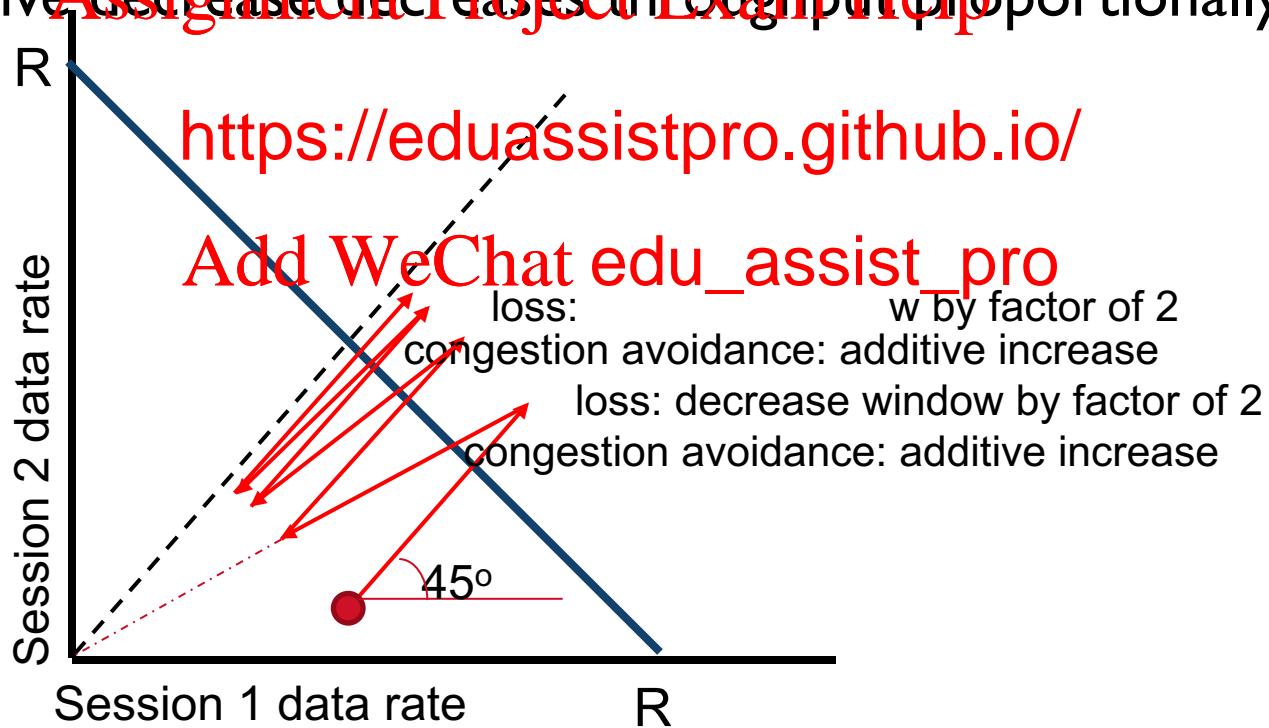
Fairness: K TCP sessions share same bottleneck link of bandwidth R, each has average rate of R/K

Assignment Project Exam Help



two competing sessions:

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



Fairness and UDP

- › multimedia apps often do not use TCP

- do not want rate

- 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

Assignment Project Exam Help

<https://eduassistpro.github.io/>

I TCP, gets 0.1R
9 TCPs, gets 0.9R

Add WeChat edu_assist_pro



Window size = min (rwnd, cwnd)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

receiveAddWindowWeChat edu_assist_pro

flow control

congestion control

- › principles behind transport layer services:
 - multiplexing, demultiplexing
 - reliable data transfer
 - connection setup, teardown
 - flow control <https://eduassistpro.github.io/>
 - congestion control [Add WeChat edu_assist_pro](#)
- › instantiation, implementation in net
 - UDP
 - TCP