

COMP 3331/9331:
Computer Networks and
Assignment Project Exam Help
Applications
WeChat: cstutorcs
Week 5

Transport Layer (Continued)

Reading Guide: Chapter 3, Sections: 3.5 – 3.7

Transport Layer 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

Practical Reliability Questions

- ❖ How do the sender and receiver keep track of outstanding pipelined segments?
- ❖ How many segments should be pipelined?
*Assignment Project Exam Help
<https://tutorcs.com>*
- ❖ How do we choose sequence numbers?
- ❖ What does connection establishment and teardown look like?
WeChat: cstutorcs
- ❖ How should we choose timeout values?

TCP: Overview

RFCs: 793, 1122, 1323, 2018, 2581

❖ **point-to-point:**

- one sender, one receiver

❖ **reliable, in-order byte stream:**

- no “message boundaries”

❖ **pipelined:**

- TCP congestion and flow control set window size

❖ **send and receive buffers**

❖ **full duplex data:**

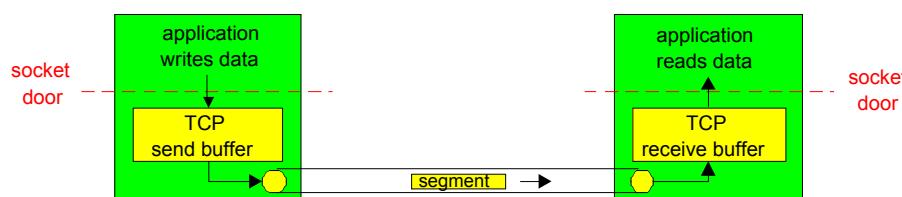
- bi-directional data flow in same connection
- MSS: maximum segment size

❖ **connection-oriented:**

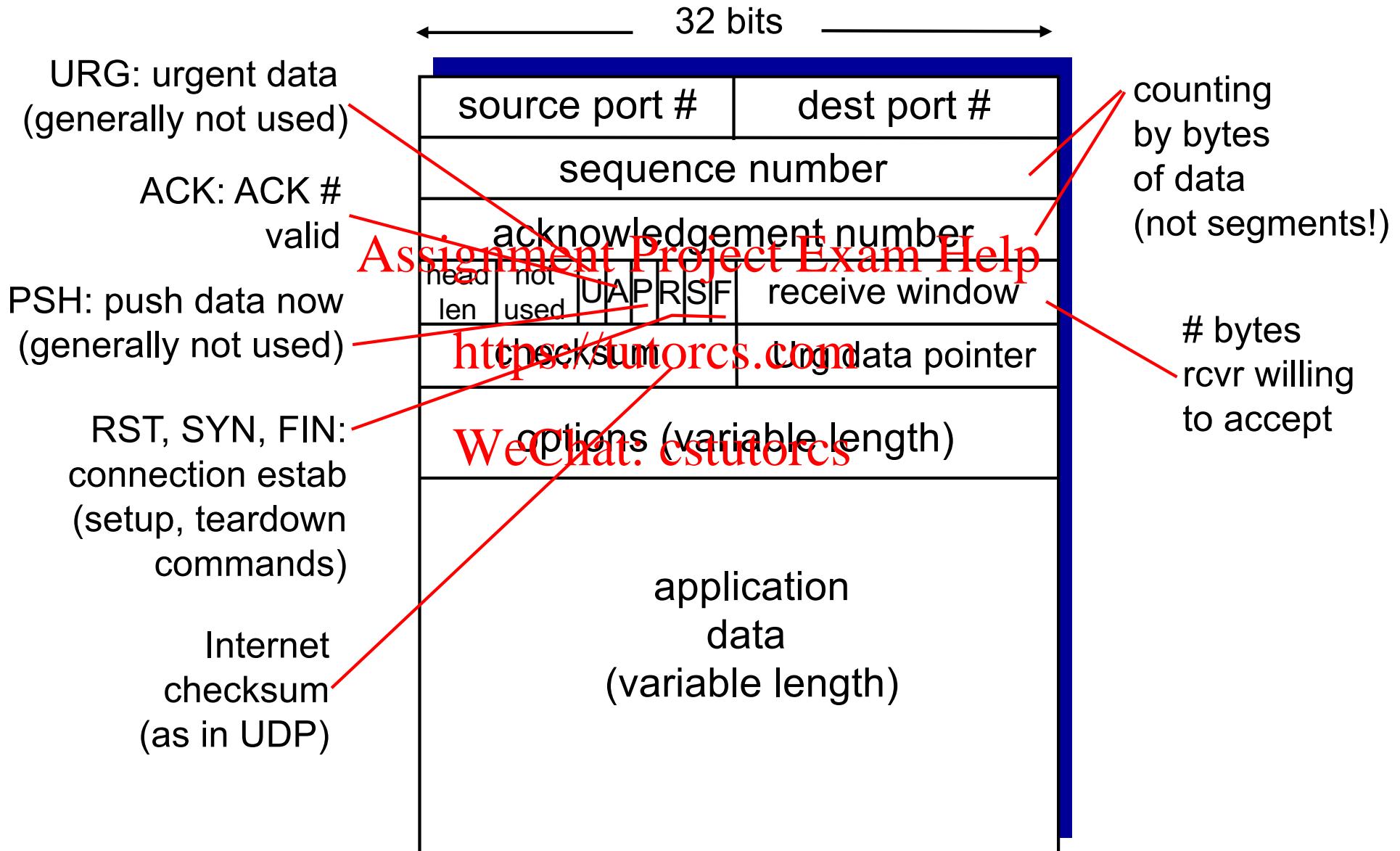
- handshaking (exchange of control msgs) init sender, receiver state before data exchange

❖ **flow controlled:**

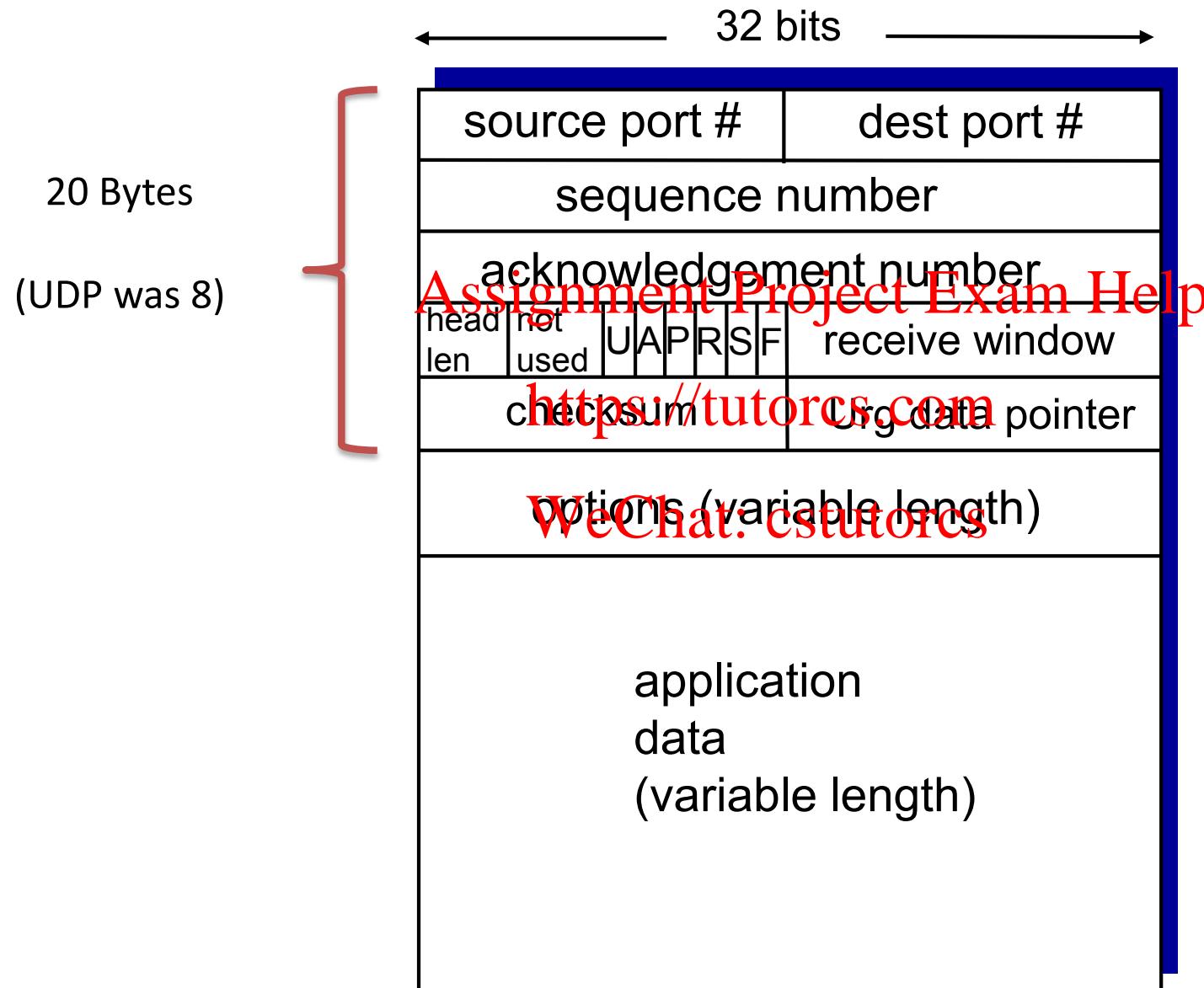
- sender will not overwhelm receiver



TCP segment structure



TCP segment structure



Transport Layer 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
- connection management

3.6 principles of congestion control

3.7 TCP congestion control

Recall: Components of a solution for reliable transport

- ❖ Checksums (for error detection)
- ❖ Timers (for loss detection)
Assignment Project Exam Help
- ❖ Acknowledgments
<https://tutorcs.com>
 - Cumulative
 - Selective WeChat: cstutorcs
- ❖ Sequence numbers (duplicates, windows)
- ❖ Sliding Windows (for efficiency)
 - Go-Back-N (GBN)
 - Selective Repeat (SR)

What does TCP do?

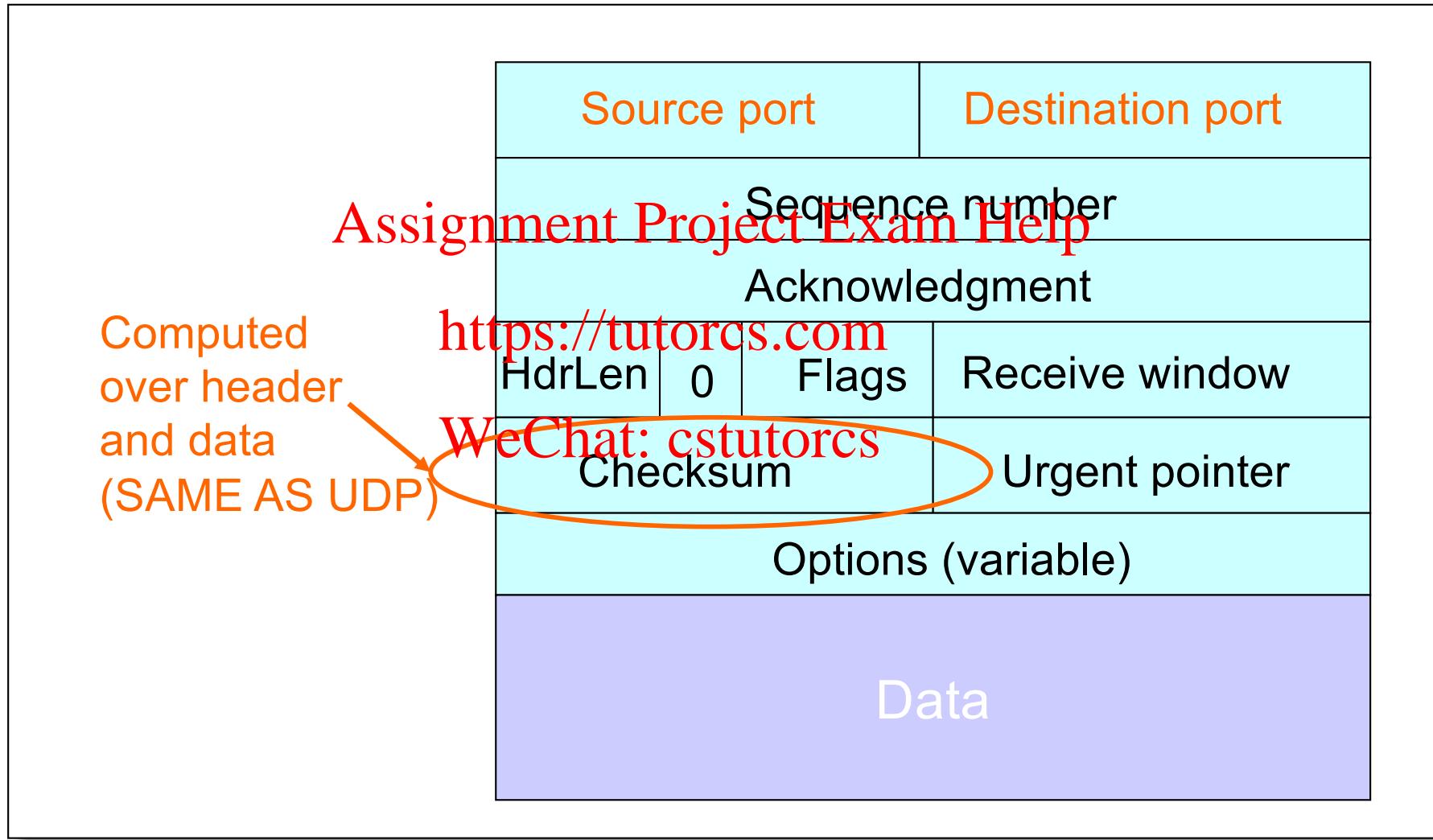
Many of our previous ideas, but some key differences

❖ Checksum Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

TCP Header



What does TCP do?

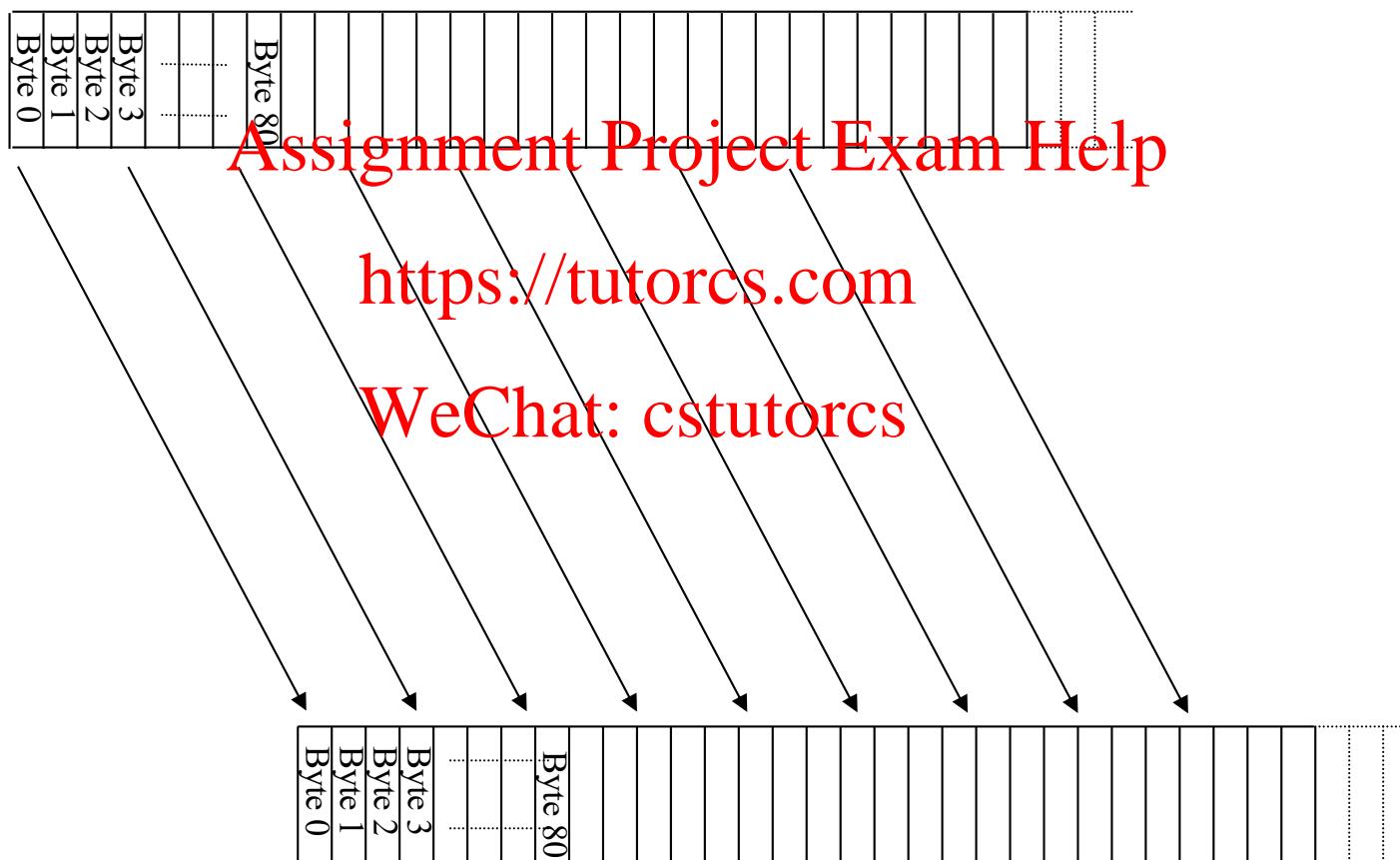
Many of our previous ideas, but some key differences

- ❖ Checksum [Assignment Project Exam Help](#)
- ❖ **Sequence numbers are byte offsets** <https://tutorcs.com>

WeChat: cstutorcs

TCP “Stream of Bytes” Service ..

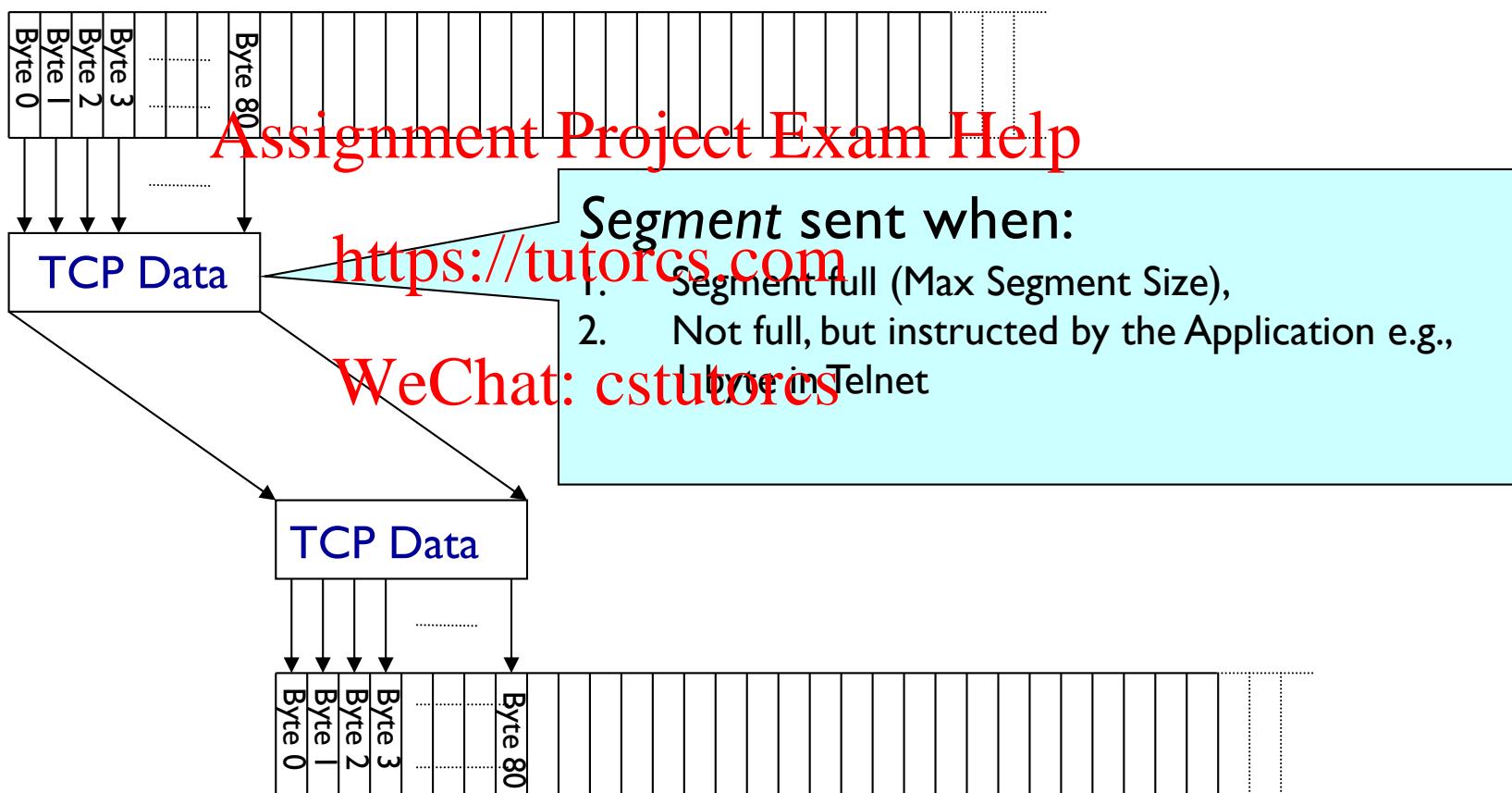
Application @ Host A



Application @ Host B

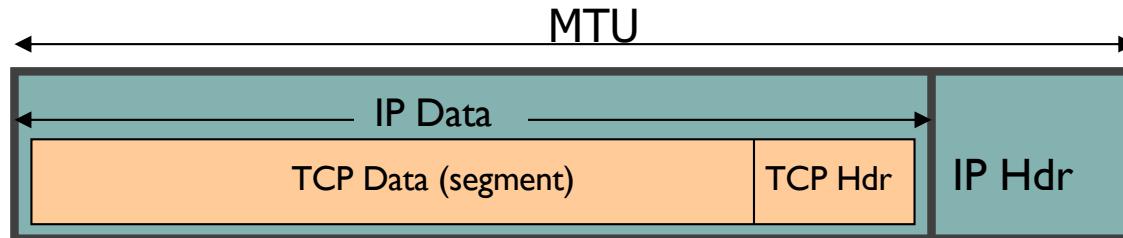
.. Provided Using TCP “Segments”

Host A



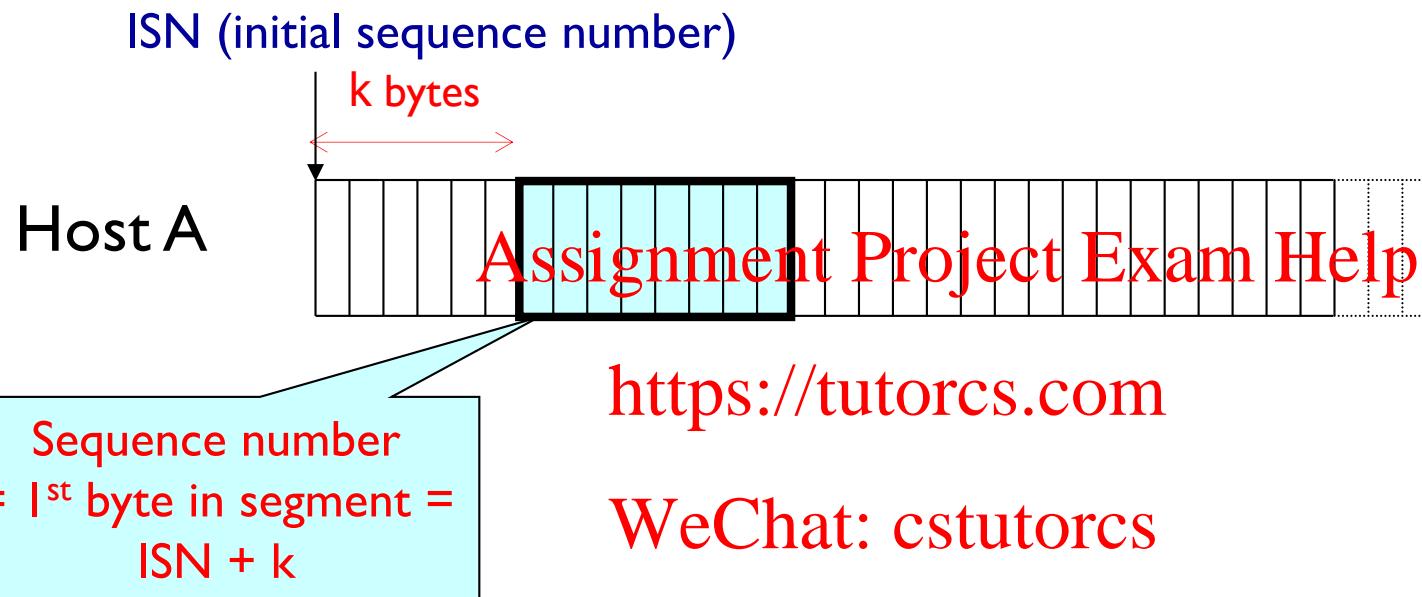
Host B

TCP Maximum Segment Size



- ❖ IP packet
 - No bigger than Maximum Transmission Unit (MTU)
 - E.g., up to 1500 bytes with Ethernet
<https://tutorcs.com>
- ❖ TCP packet
 - IP packet with a TCP header and data inside
 - TCP header \geq 20 bytes long
- ❖ TCP segment
 - No more than Maximum Segment Size (MSS) bytes
 - E.g., up to 1460 consecutive bytes from the stream
 - $MSS = MTU - 20 \text{ (min IP header)} - 20 \text{ (min TCP header)}$

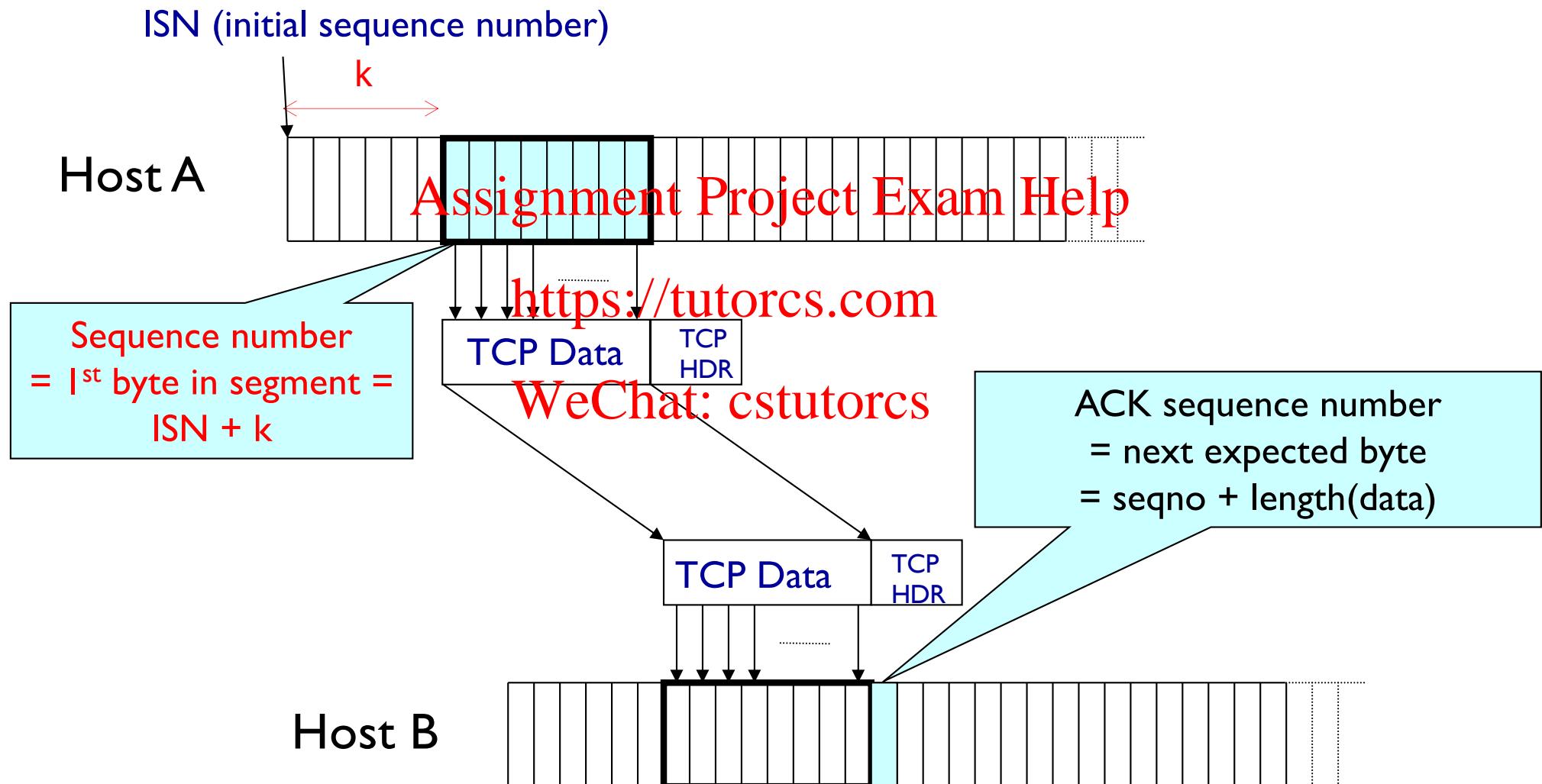
Sequence Numbers



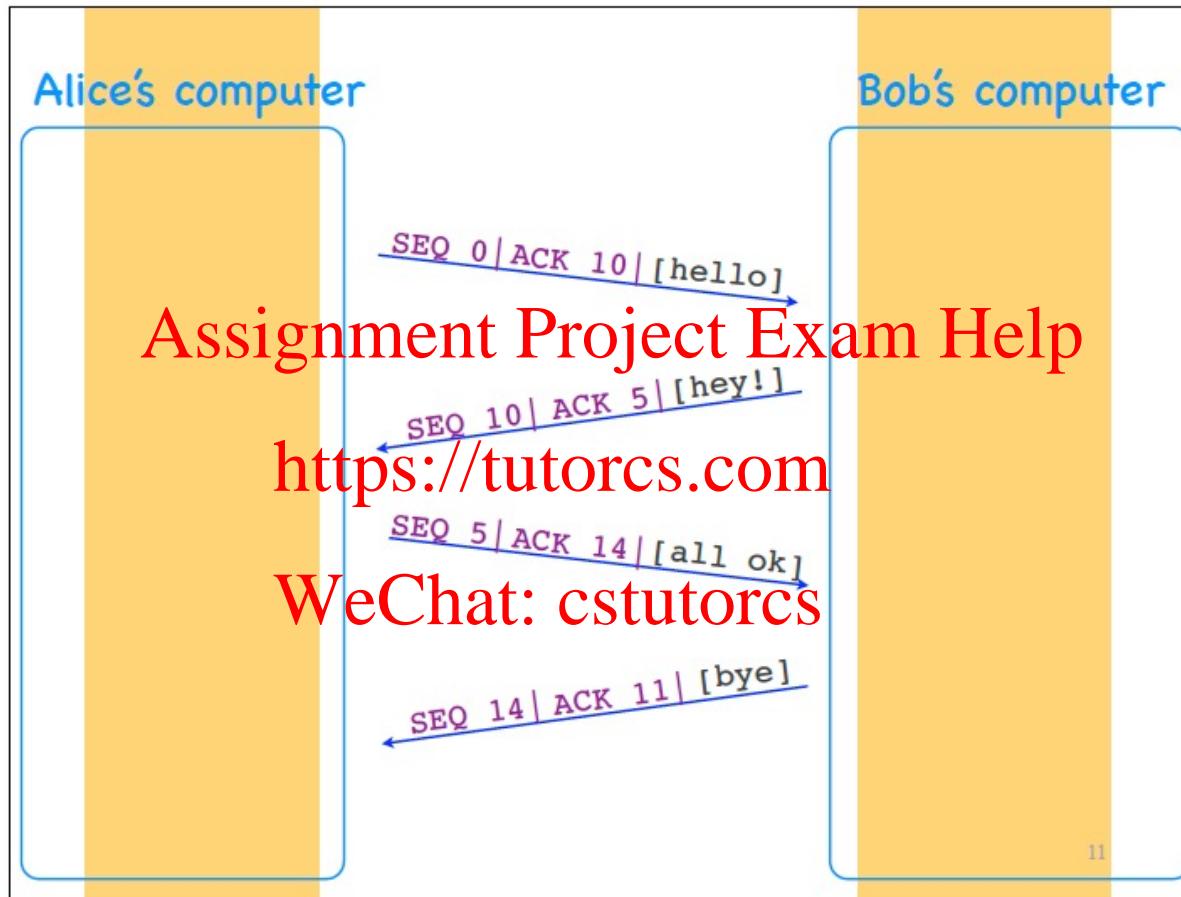
Sequence numbers:

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

Sequence & Ack Numbers

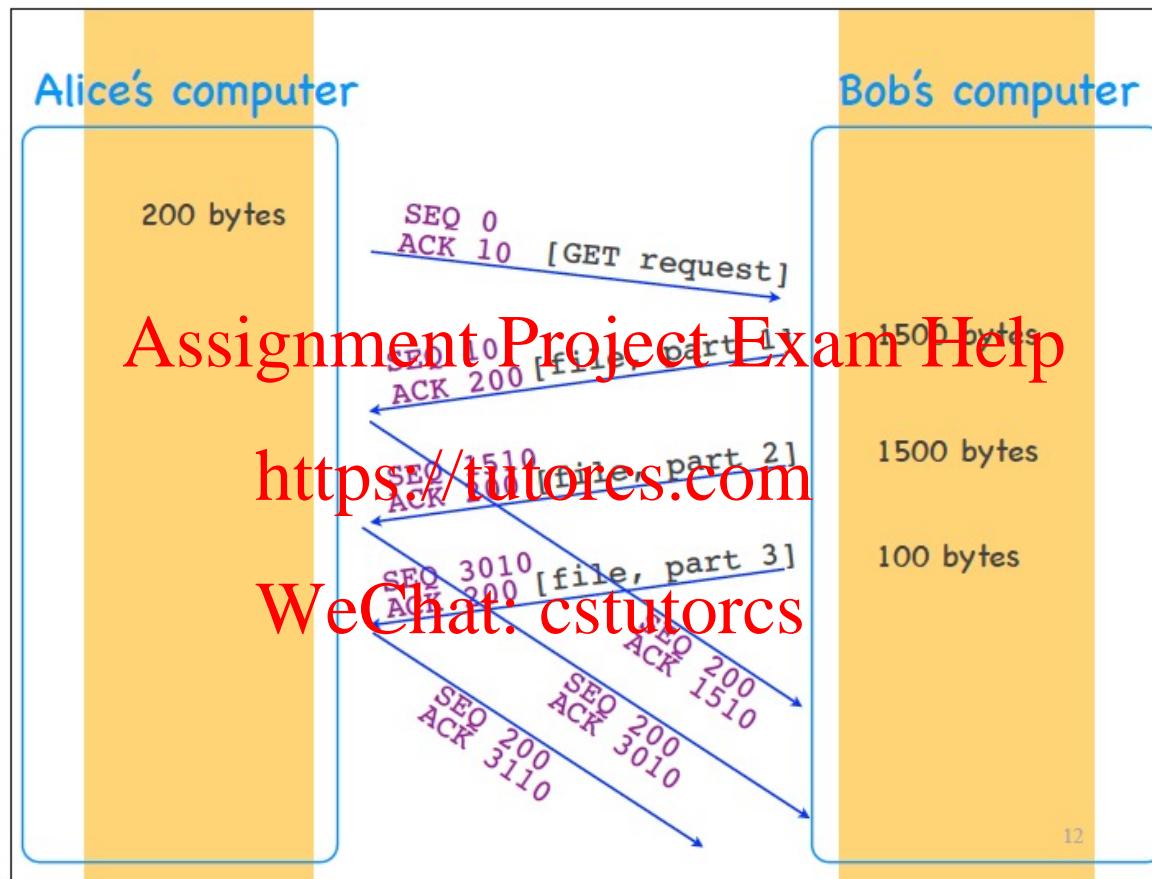


Example



Note: Connection establishment not shown. Alice's end point selects the initial sequence number as 0 while Bob's end point selects the initial sequence number as 10

Another Example

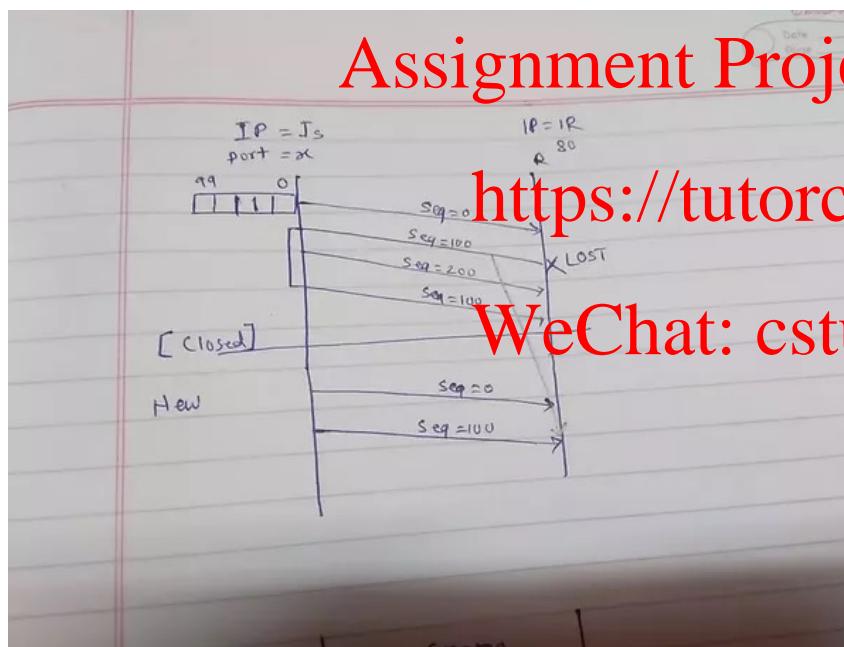


Note: Connection establishment not shown. Alice's end point selects the initial sequence number as 0 while Bob's end point selects the initial sequence number as 10

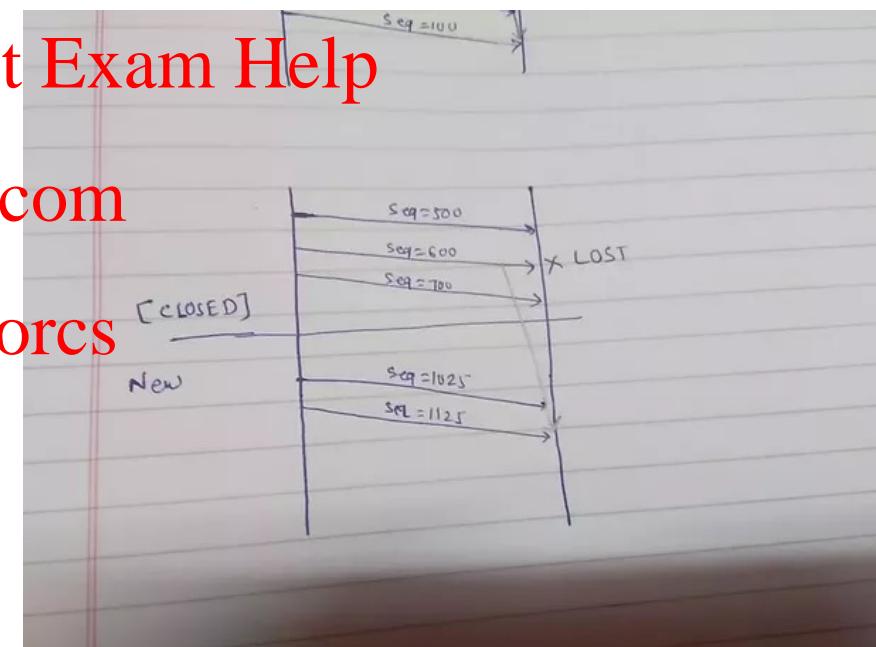
HTTP response split into 3 segments (MSS = 1500 bytes)

Why choose random ISN?

- ❖ Avoids ambiguity with back-to-back connections between same end-points



(a) When ISN=0



(b) When ISN is random

- ❖ Potential security issue if the ISN is known

What does TCP do?

Most of our previous tricks, but a few differences

- ❖ Checksum
- ❖ Sequence numbers are byte offsets
- ❖ Receiver sends cumulative acknowledgements (like GBN)

Assignment Project Exam Help

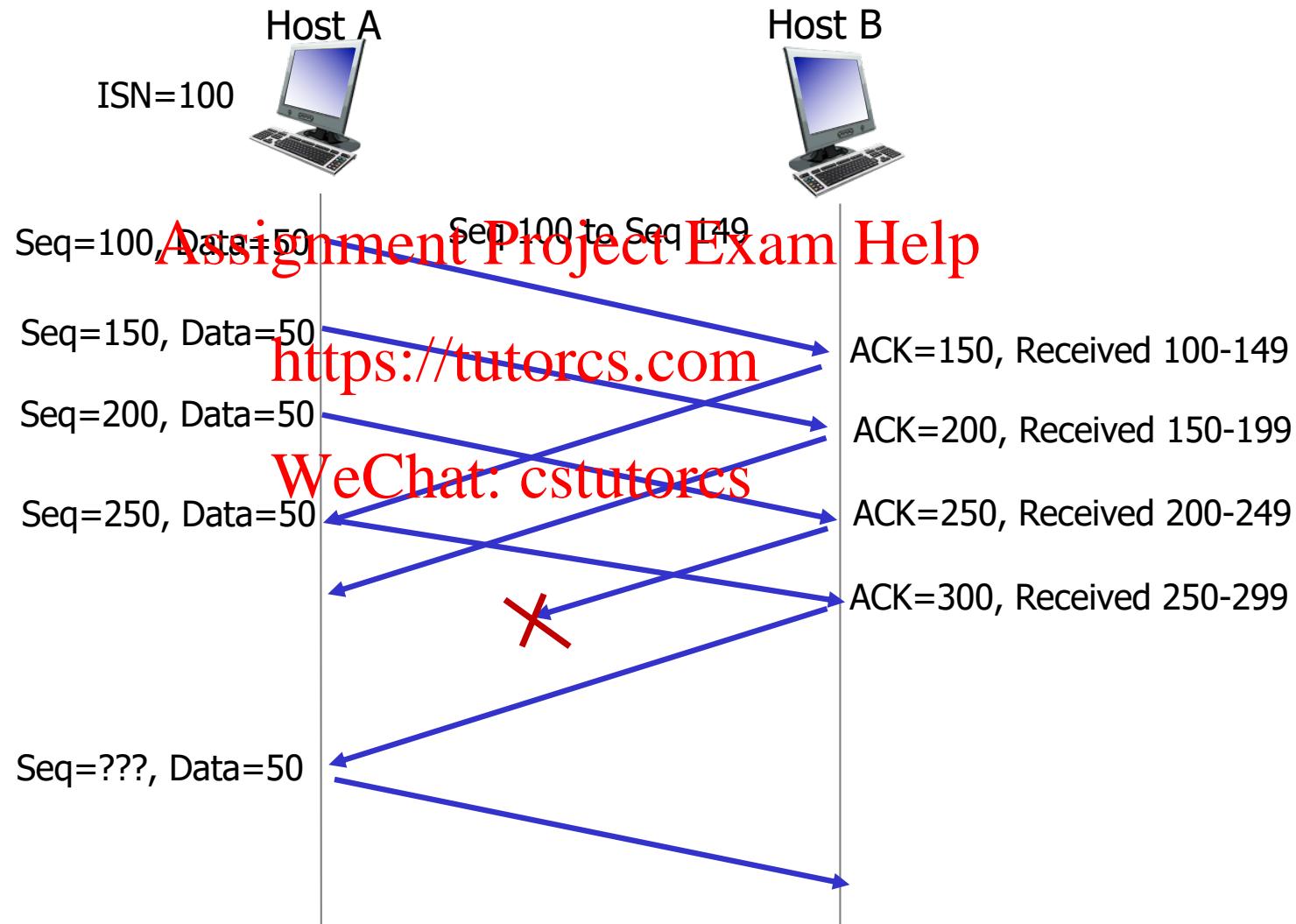
<https://tutorcs.com>

WeChat: cstutorcs

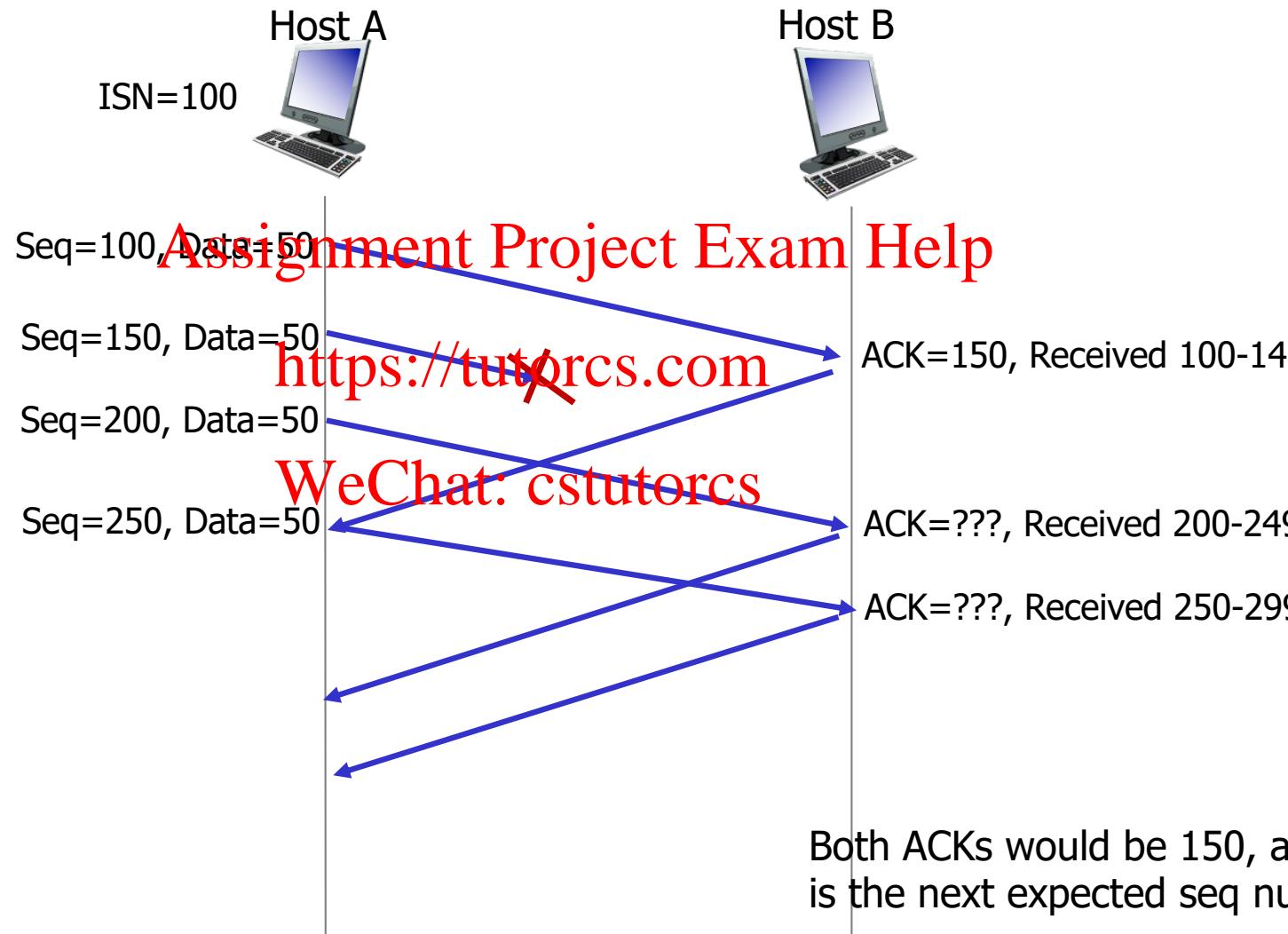
ACKing and Sequence Numbers

- ❖ Sender sends packet
 - Data starts with sequence number X
 - Packet contains B bytes [X, X+1, X+2, ..., X+B-1]
- Assignment Project Exam Help
- ❖ Upon receipt of packet, receiver sends an ACK
 - If all data prior to X already received.
 - ACK acknowledges X+B (because that is next expected byte)
 - If highest in-order byte received is Y s.t. (Y+1) < X
 - ACK acknowledges Y+1
 - Even if this has been ACKed before

TCP seq. numbers, ACKs



TCP seq. numbers, ACKs



Normal Pattern

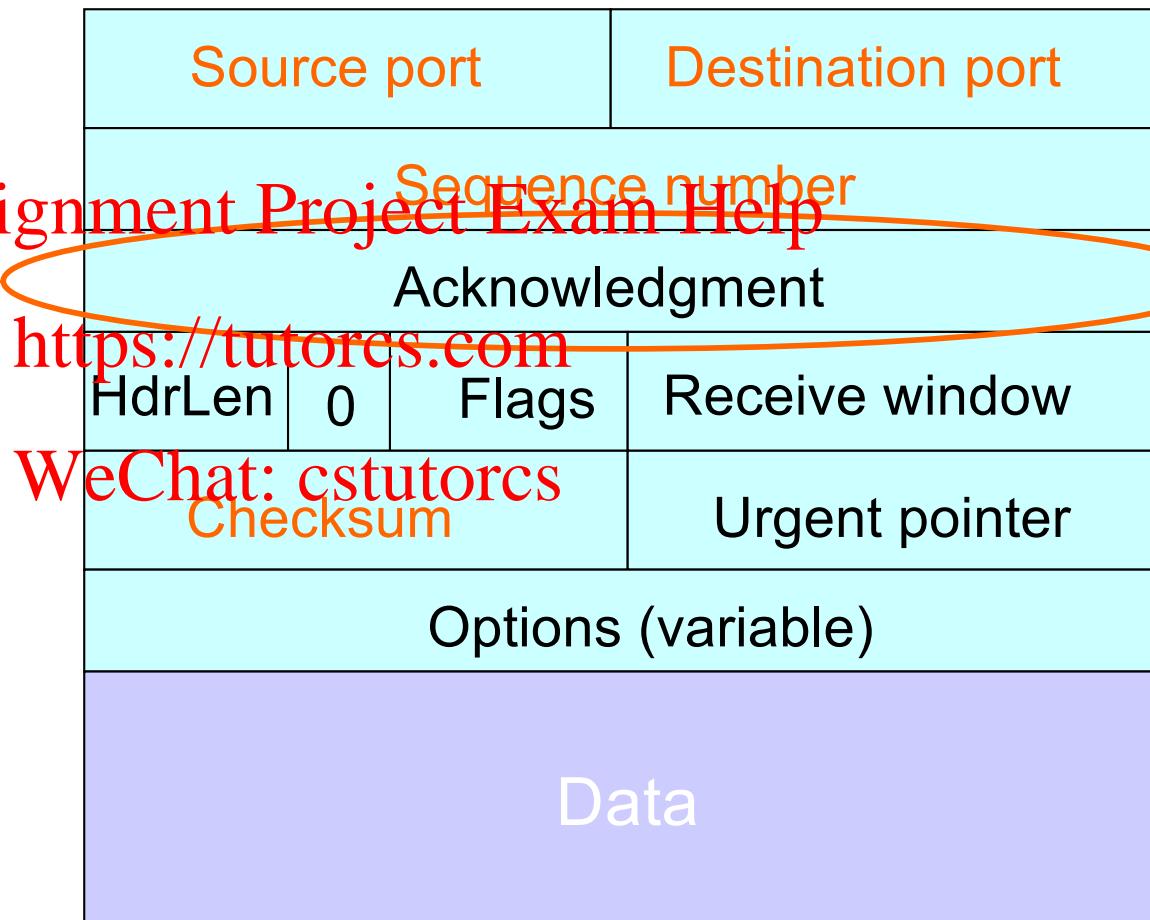
- ❖ Sender: seqno=X, length=B
- ❖ Receiver: ACK=X+B
- ❖ Sender: seqno=X+B, length=B
- ❖ Receiver: ACK=X+2B
- ❖ Sender: seqno=X+2B, length=B
- ❖ Seqno of next packet is same as last ACK field

Packet Loss

- ❖ Sender: seqno=X, length=B
- ❖ Receiver: ACK=X+B
- ❖ Sender: ~~seqno=X+B, length=B~~ LOST
Assignment Project Exam Help
- ❖ Sender: seqno=X+2B, length=B
<https://tutorcs.com>
- ❖ Receiver: ACK = X+B
WeChat: cstutorcs

TCP Header

Acknowledgment gives seqno just beyond highest seqno received **in order** ("What Byte is Next")



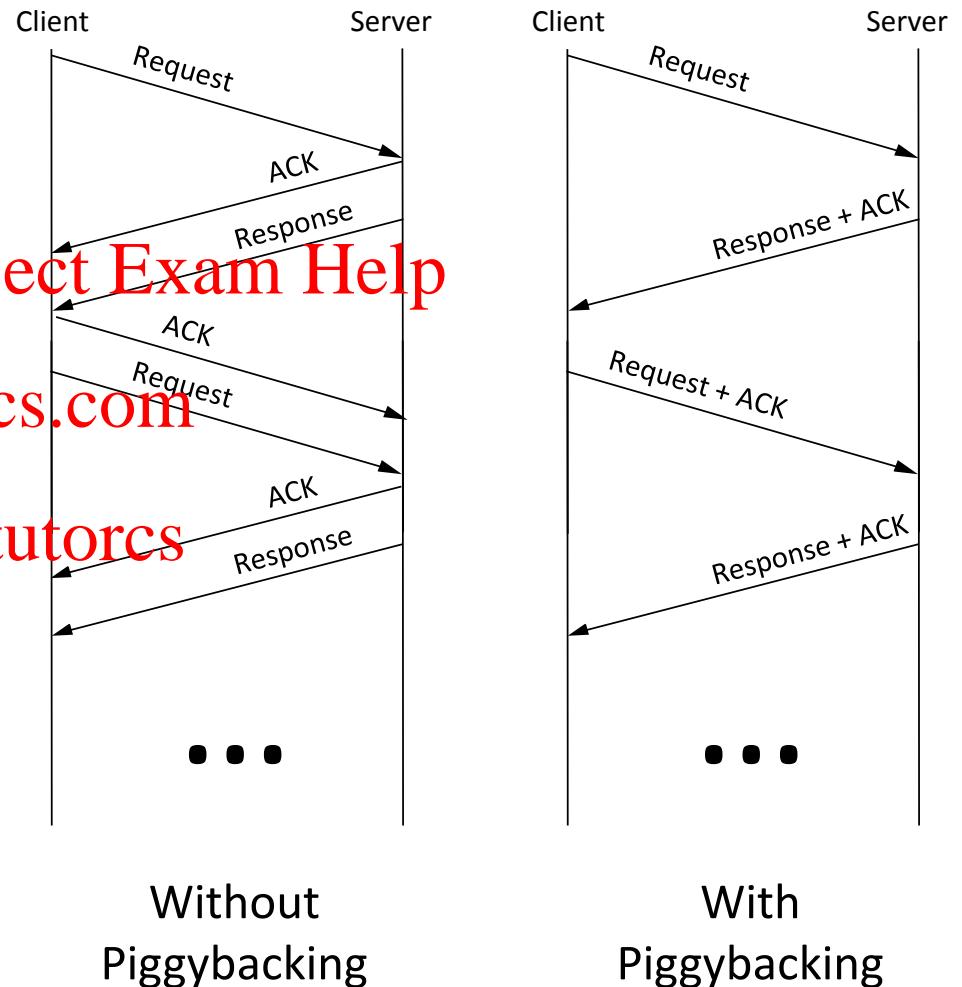
Piggybacking

- ❖ So far, we've assumed distinct “sender” and “receiver” roles

Assignment Project Exam Help

- ❖ In reality, usually both sides of a connection send some data

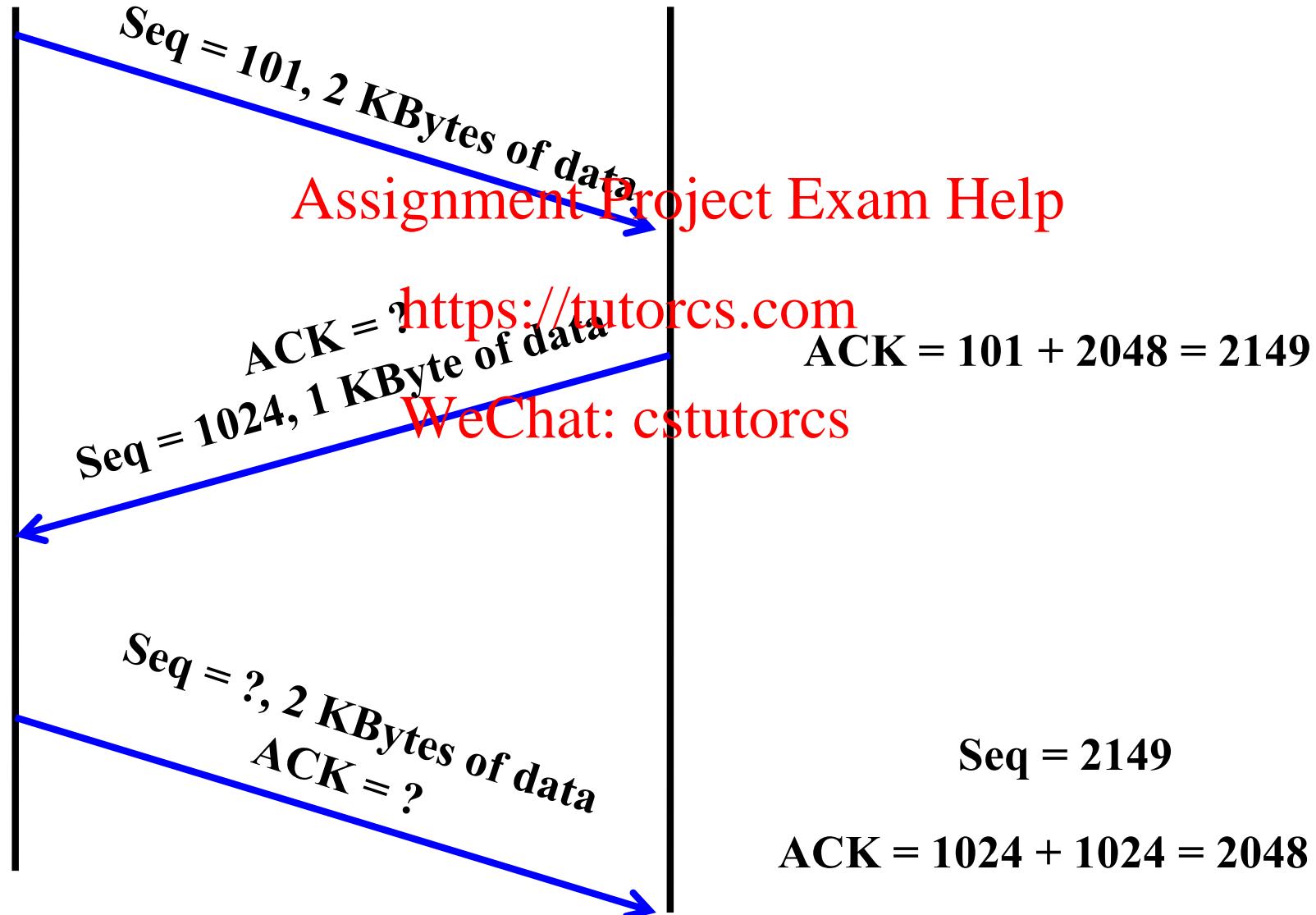
<https://tutorcs.com>
WeChat: cstutorcs



Without
Piggybacking

With
Piggybacking

Quiz



What does TCP do?

Most of our previous tricks, but a few differences

- ❖ Checksum
- ❖ Sequence numbers are byte offsets
- ❖ Receiver sends cumulative acknowledgements (like GBN)
- ❖ Receivers can buffer out-of-sequence packets (like SR)

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

Loss with cumulative ACKs

- ❖ Sender sends packets with 100Bytes and sequence numbers:
 - 100, 200, 300, 400, 500, 600, 700, 800, 900, ...
<https://tutorcs.com>
- ❖ Assume the ~~first~~ packet (seq. no. 500) is lost, but no others
- ❖ Stream of ACKs will be:
 - 200, 300, 400, 500, 500, 500, ...

What does TCP do?

Most of our previous tricks, but a few differences

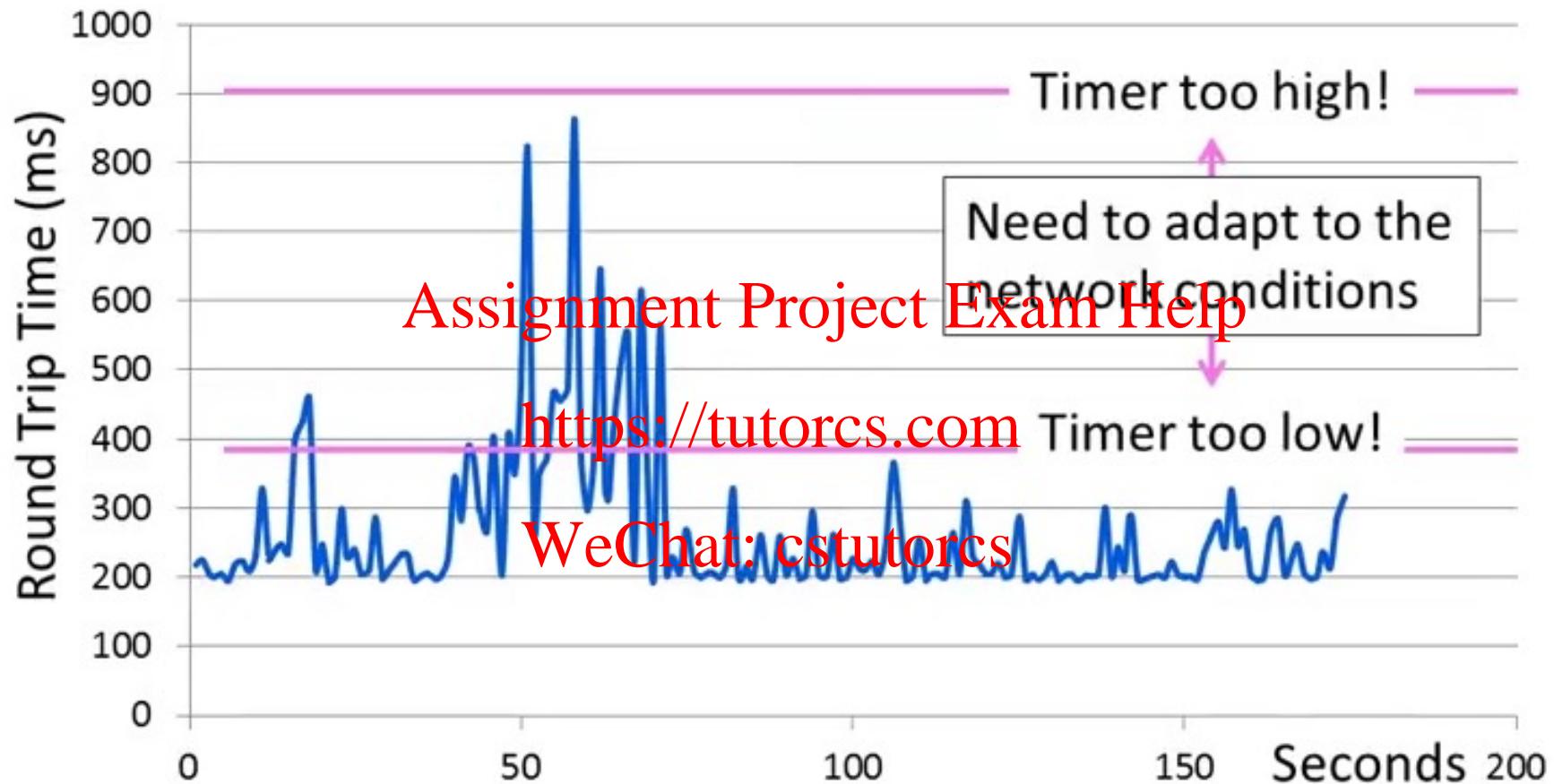
- ❖ Checksum
- ❖ Sequence numbers are byte offsets
- ❖ Receiver sends cumulative acknowledgements (like GBN)
- ❖ Receivers do not drop out-of-sequence packets (like SR)
- ❖ Sender maintains a single retransmission timer (like GBN) and retransmits on timeout (*how much?*)

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

TCP round trip time, timeout



TCP round trip time, timeout

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 and connection has lower throughput

Q: how to estimate RTT?

- ❖ **SampleRTT**: measured time from segment transmission until ACK receipt

<https://tutorcs.com> ignore retransmissions

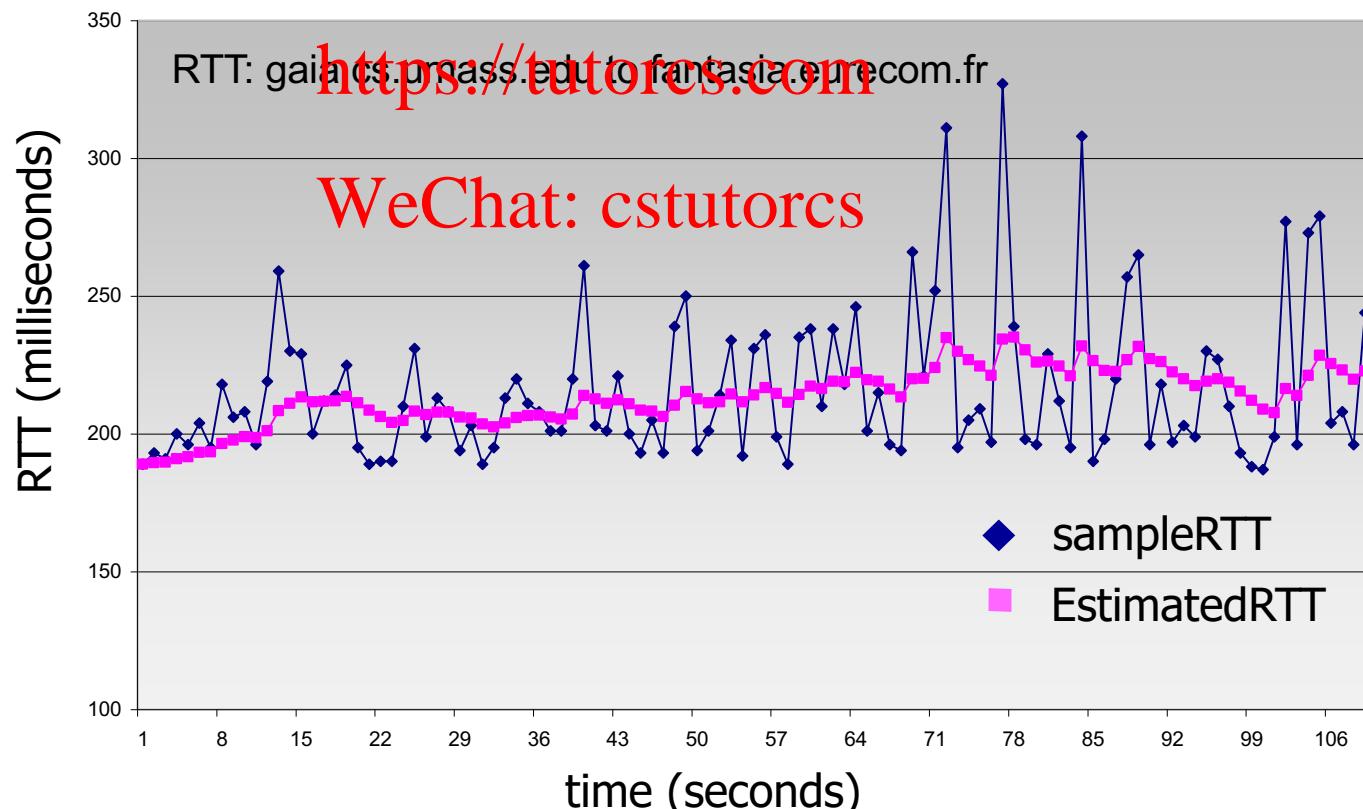
WeChat: cstutorcs
❖ **SampleRTT** will vary, want estimated RTT “smoother”

- average several *recent* measurements, not just current **SampleRTT**

TCP round trip time, timeout

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

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



TCP round trip time, timeout

- ❖ **timeout interval:** **EstimatedRTT** plus “safety margin”
 - large variation in **EstimatedRTT** -> larger safety margin

- ❖ estimate SampleRTT deviation from EstimatedRTT:
Assignment Project Exam Help

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

(typically, $\beta = 0.25$)

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



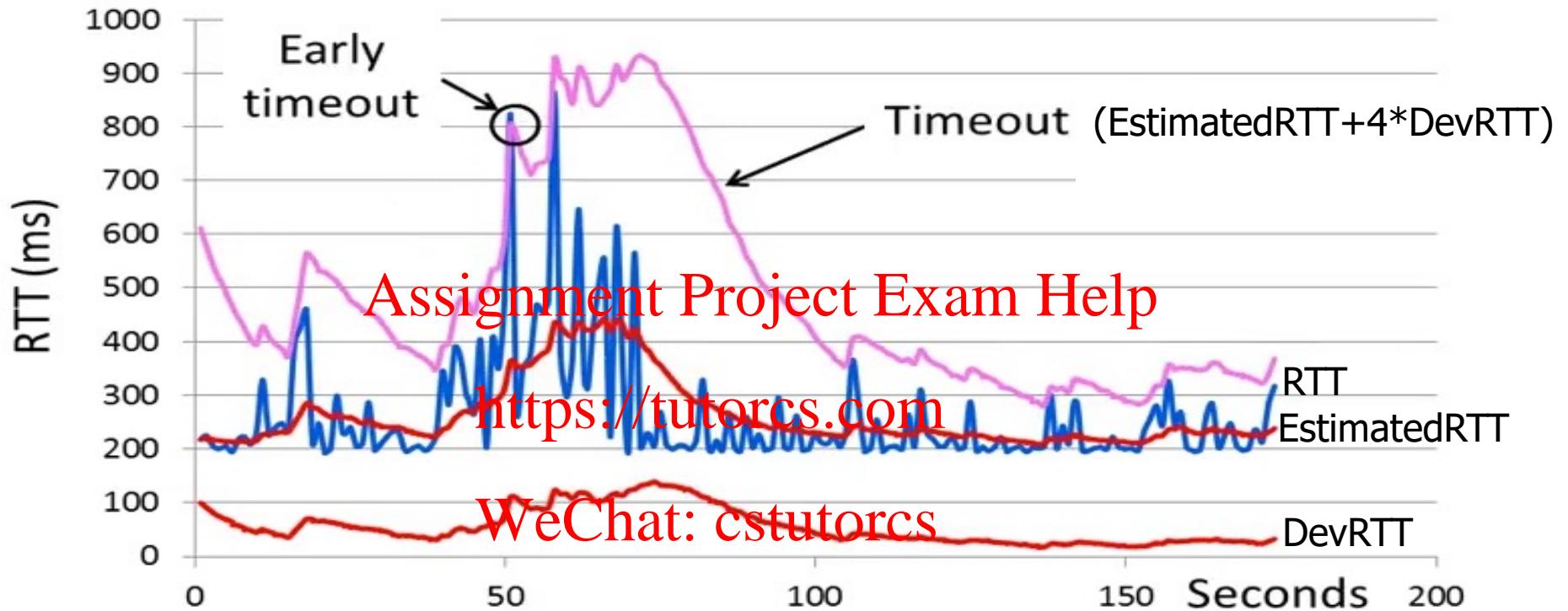
↑
estimated RTT

↑
“safety margin”

Practice Problem:

http://wps.pearsoned.com/ecs_kurose_compnetw_6/216/55463/14198700.cw/index.html

TCP round trip time, timeout



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



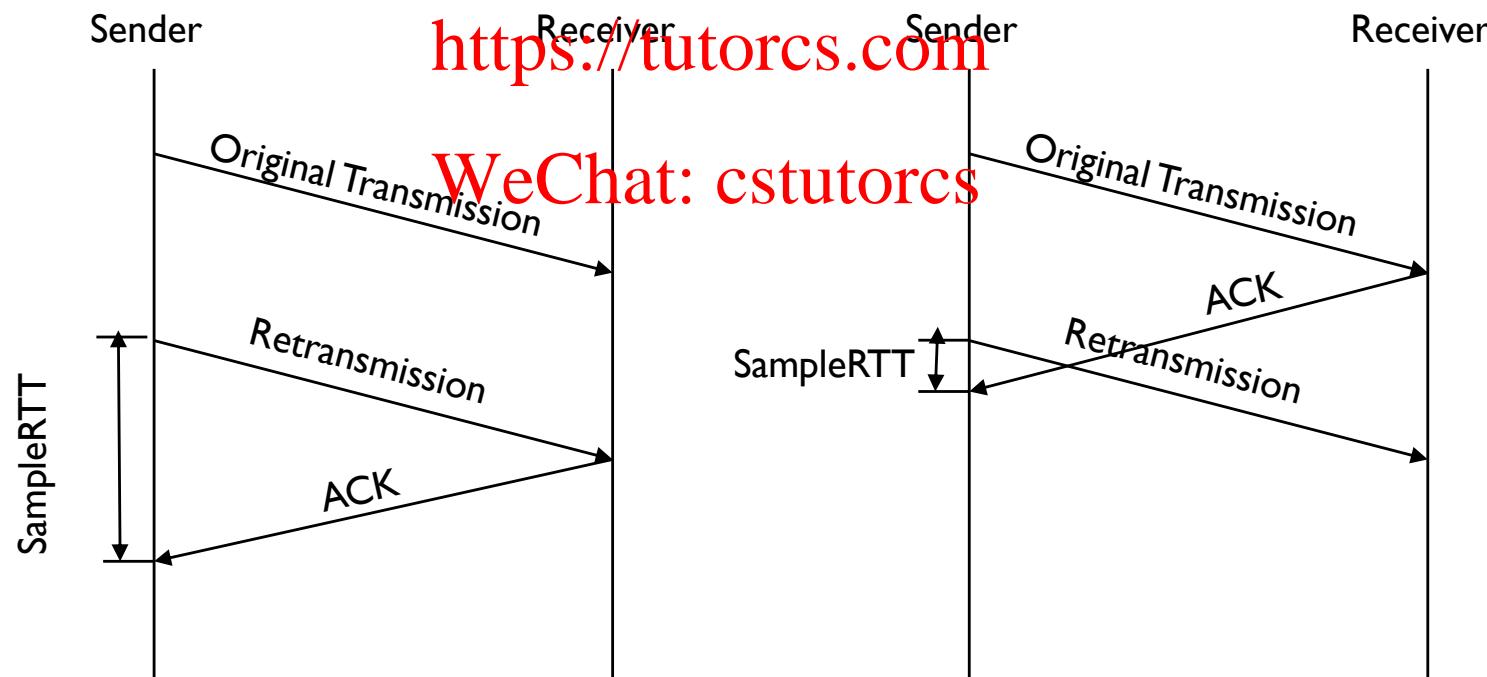
estimated RTT

“safety margin”

Why exclude retransmissions in RTT computation?

- ❖ How do we differentiate between the real ACK, and ACK of the retransmitted packet?

Assignment Project Exam Help



TCP sender events:

PUTTING IT
TOGETHER

data rcvd from app:

- ❖ 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

timeout:

- ❖ retransmit segment that caused timeout

restart timer

ack rcvd:

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

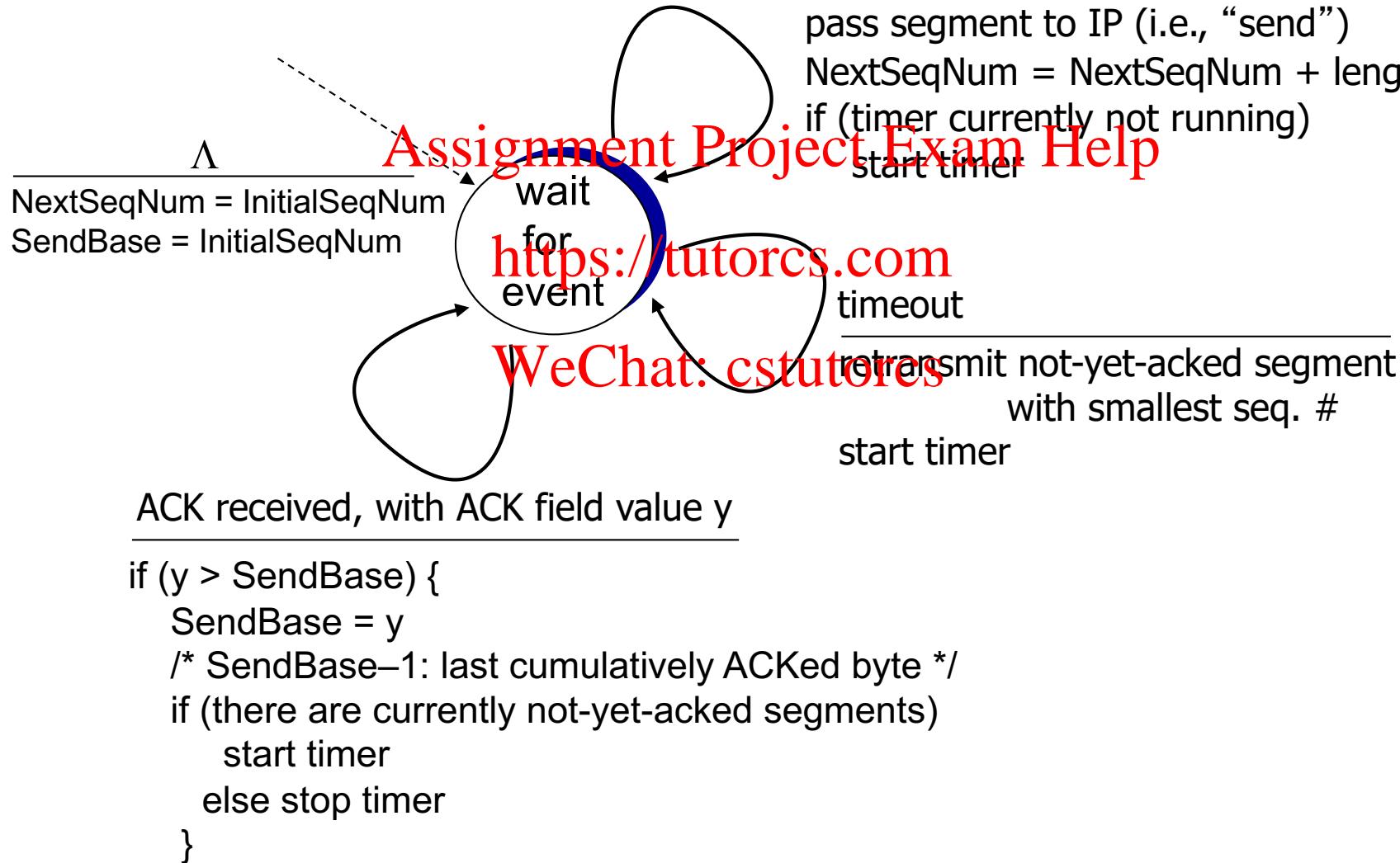
Assignment Project Exam Help

<https://tutorcs.com>

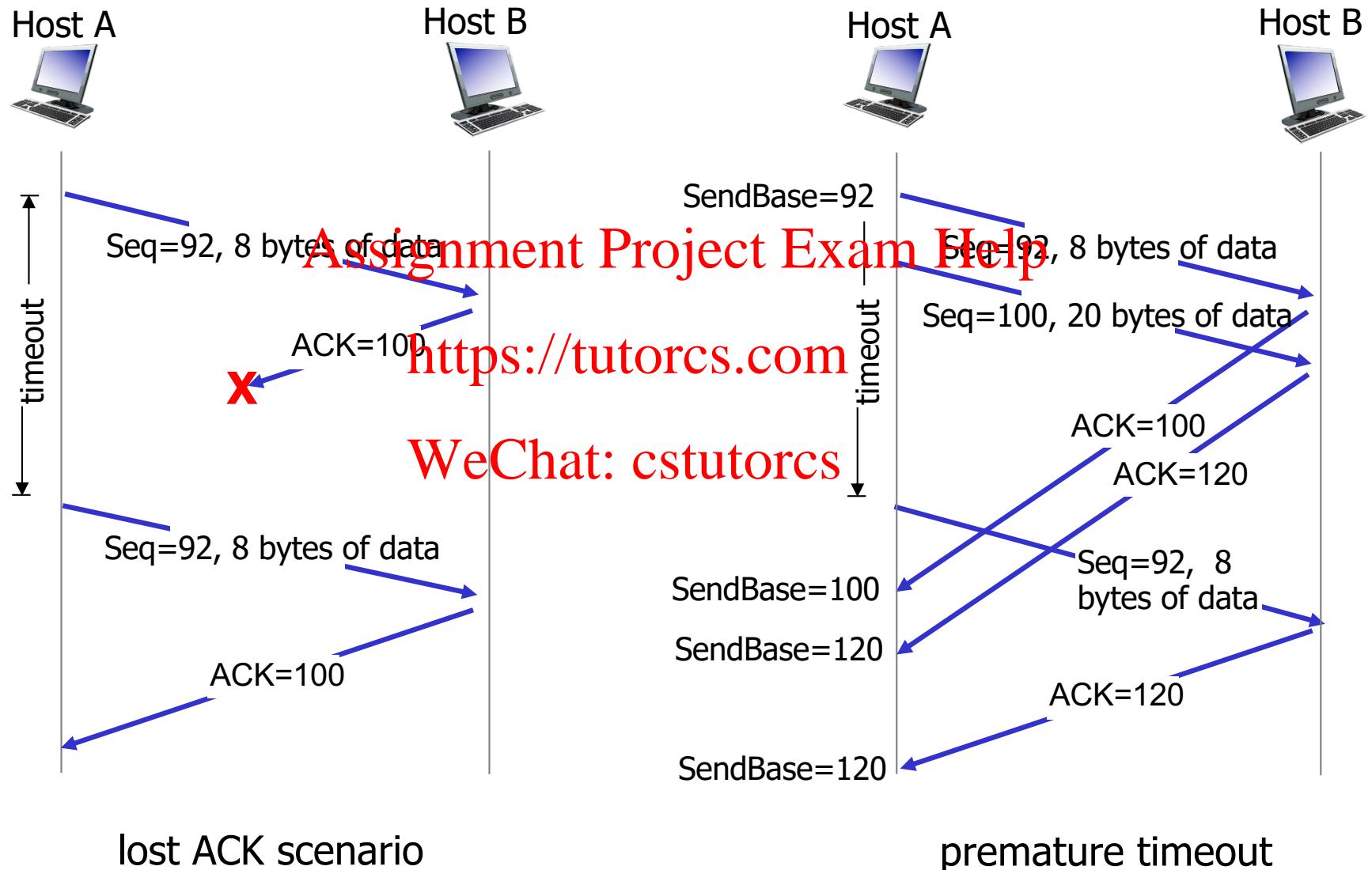
WeChat: cstutores

TCP sender (simplified)

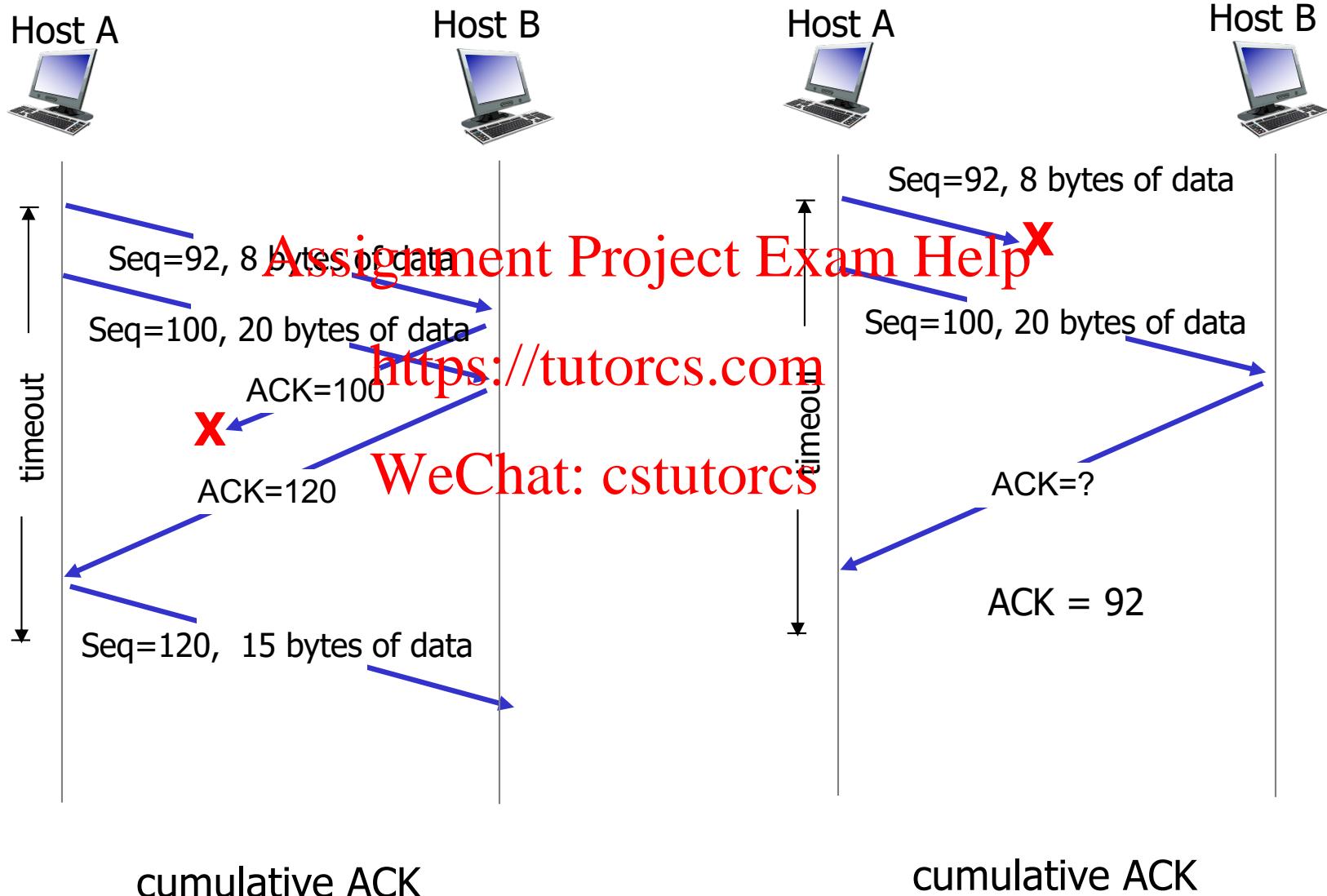
PUTTING IT
TOGETHER



TCP: retransmission scenarios



TCP: retransmission scenarios



TCP ACK generation [RFC 1122, RFC 2581]

<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	immediately send single cumulative ACK , ACKing both in-order segments
arrival of out-of-order segment higher-than-expect seq. # . Gap detected	immediately send duplicate ACK , 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

What does TCP do?

Most of our previous tricks, but a few differences

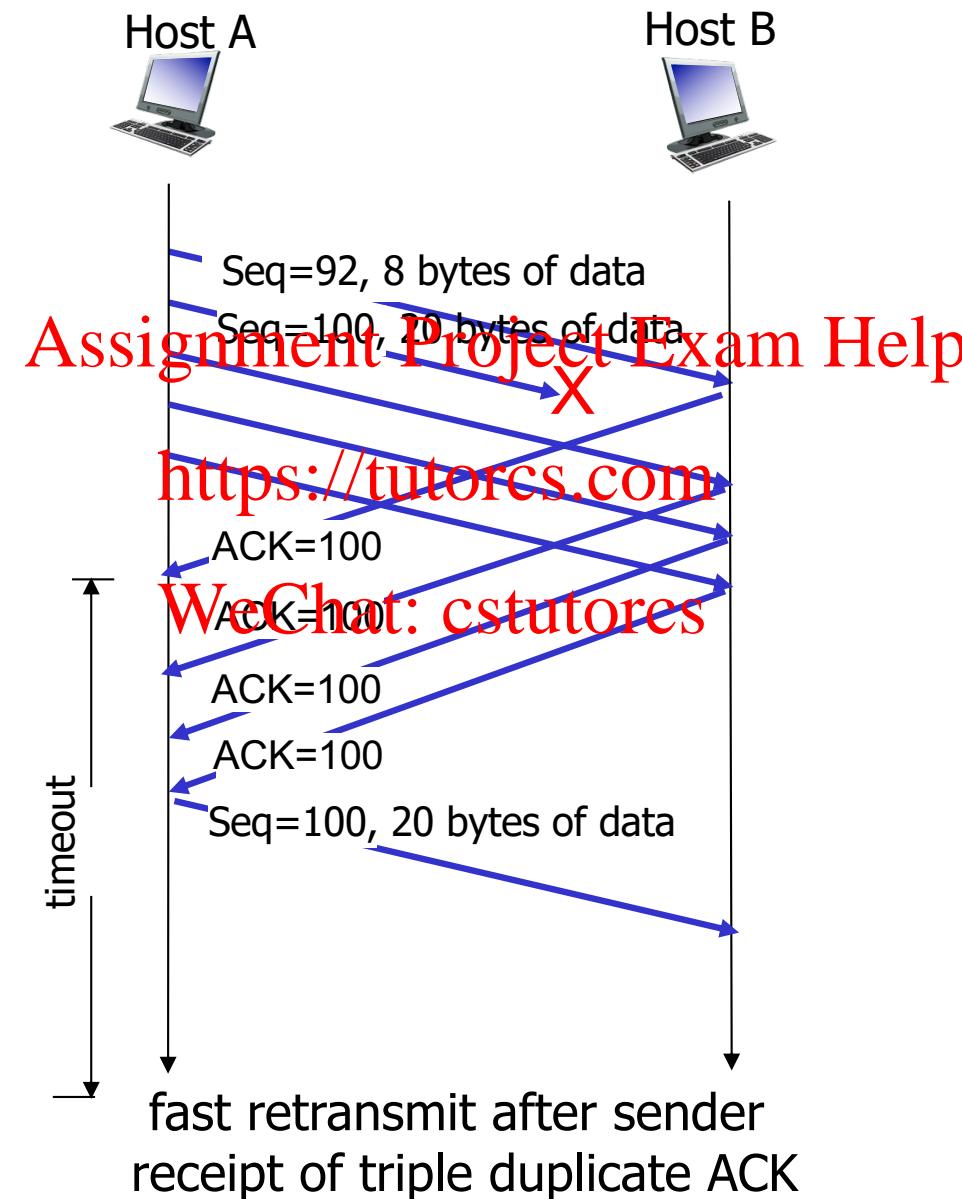
- ❖ Checksum
- ❖ Sequence numbers are byte offsets
- ❖ Receiver sends cumulative acknowledgements (like GBN)
- ❖ Receivers may not drop out-of-sequence packets (like SR)
- ❖ Sender maintains a single retransmission timer (like GBN) and retransmits on timeout
- ❖ Introduces **fast retransmit**: optimisation that uses duplicate ACKs to trigger early retransmission

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

TCP fast retransmit



TCP fast retransmit

- ❖ time-out period often relatively long:
 - long delay before resending lost packet
- ❖ “Duplicate ACKs” are a sign of an isolated loss
 - The lack of ACK <https://tutorcs.com> means that packet hasn’t been delivered
 - Stream of ACKs means some packets are being delivered
 - Could trigger resend on receiving “k” duplicate ACKs (TCP uses k = 3)

TCP fast retransmit

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

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

What does TCP do?

Most of our previous ideas, but some key differences

- ❖ Checksum [Assignment Project Exam Help](https://tutorcs.com)
- ❖ Sequence numbers <https://tutorcs.com>
- ❖ Receiver sends cumulative acknowledgements (like GBN)
- ❖ Receivers do not drop out-of-sequence packets (like SR)
- ❖ Sender maintains a single retransmission timer (like GBN) and retransmits on timeout
- ❖ Introduces fast retransmit: optimization that uses duplicate ACKs to trigger early retransmission

Quiz: TCP Sequence Numbers?



A TCP Sender is just about to send a segment of size 100 bytes with sequence number 1234 and ack number 436 in the TCP header. What is the highest sequence number up to (and including) which this sender has received all bytes from the receiver?
<https://tutorcs.com>

- A. 1233
- B. 436
- C. 435
- D. 1334
- E. 536

WeChat: cstutorcs

ANSWER: C

Quiz: TCP Sequence Numbers?



A TCP Sender is just about to send a segment of size 100 bytes with sequence number 1234 and ack number 436 in the TCP header. Is it possible that the receiver has received byte number 1335? Assignment Project Exam Help

- A. Yes
- B. No

<https://tutorcs.com>

WeChat: cstutorcs

ANSWER: A

Quiz: TCP Sequence Numbers?



The following statement is true about the TCP sliding window protocol for implementing reliable data transfer

- A. It exclusively uses the ideas of Go-Back-N
- B. It exclusively uses the ideas of Selective Repeat
- C. It uses a combination of ideas of Go-Back-N and Selective-Repeat
- D. It uses none of the ideas of Go-Back-N and Selective-Repeat

ANSWER: C

Transport Layer 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

<https://tutorcs.com>

- flow control
- connection management

WeChat: cstutorcs

3.6 principles of congestion control

3.7 TCP congestion control

TCP flow control

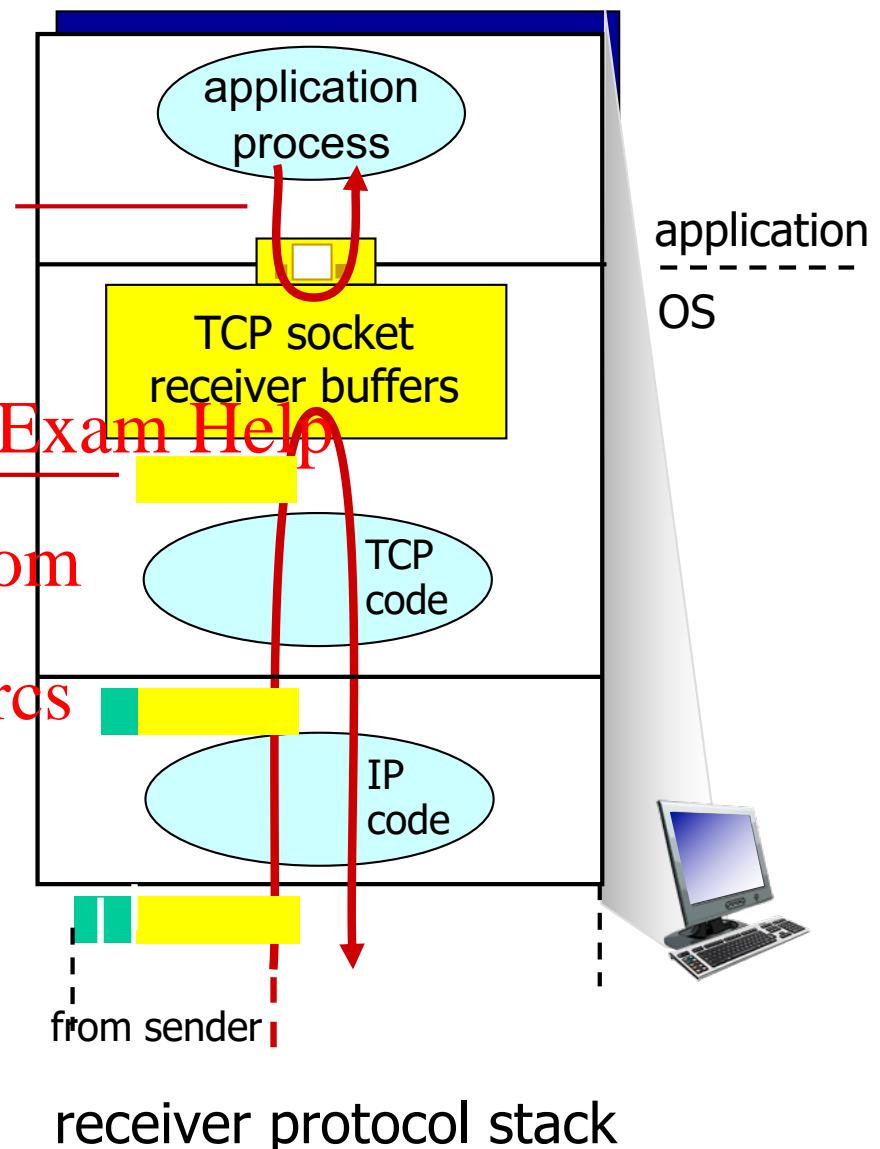
application may
remove data from
TCP socket buffers

Assignment Project Exam Help
... slower than TCP
receiver is delivering
(sender is sending)
<https://tutorcs.com>

WeChat: cstutorcs

flow control

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



TCP flow control

- ❖ receiver “advertises” free buffer space by including **rwnd** value in TCP header of receiver-to-sender segments

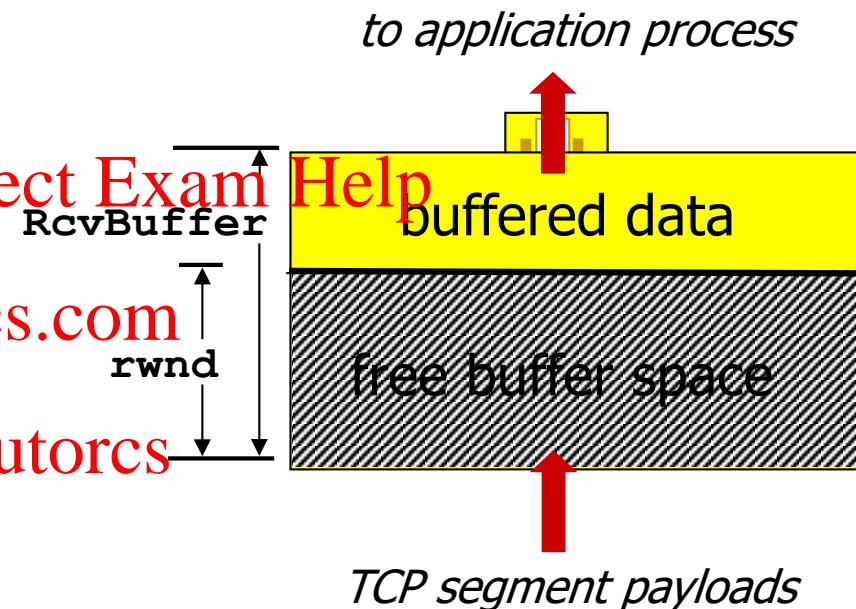
- **RcvBuffer** size set via socket options (typical default is 4096 bytes)
- many operating systems autoadjust **RcvBuffer**

- ❖ sender limits amount of unacked (“in-flight”) data to receiver’s **rwnd** value
- ❖ guarantees receive buffer will not overflow

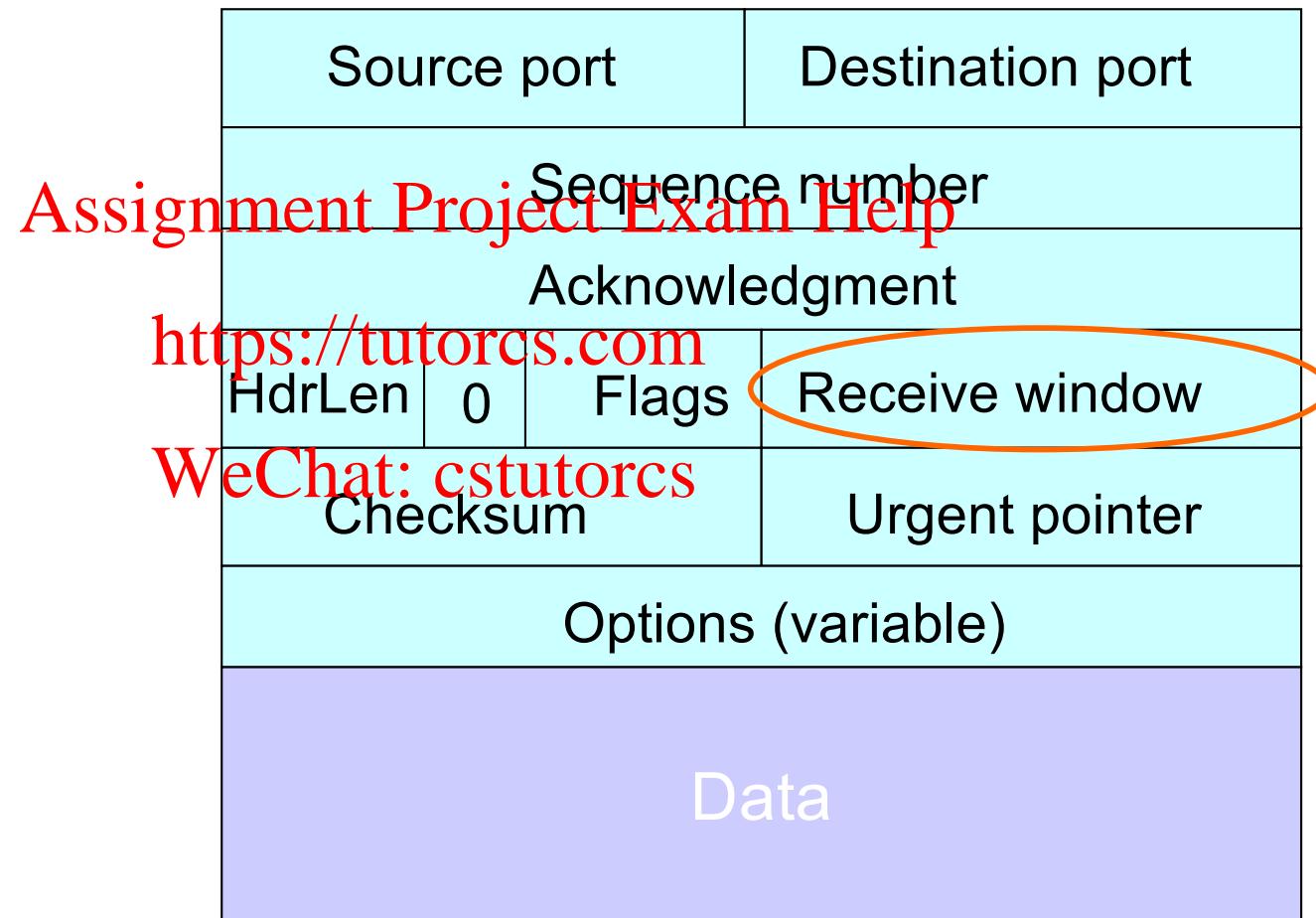
Assignment Project Exam

<https://tutorcs.com>

WeChat: cstutorcs



TCP Header



TCP flow control

- ❖ What if **rwnd** = 0?
 - Sender would stop sending data
 - Eventually the receive buffer would have space when the application process reads some bytes
 - But how does <https://tutorcs.com> advertise the new **rwnd** to the sender?
- ❖ Sender keeps sending TCP segments with one data byte to the receiver
- ❖ These segments are dropped but acknowledged by the receiver with a zero-window size
- ❖ Eventually when the buffer empties, non-zero window is advertised

Transport Layer 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

Assignment Project Exam Help

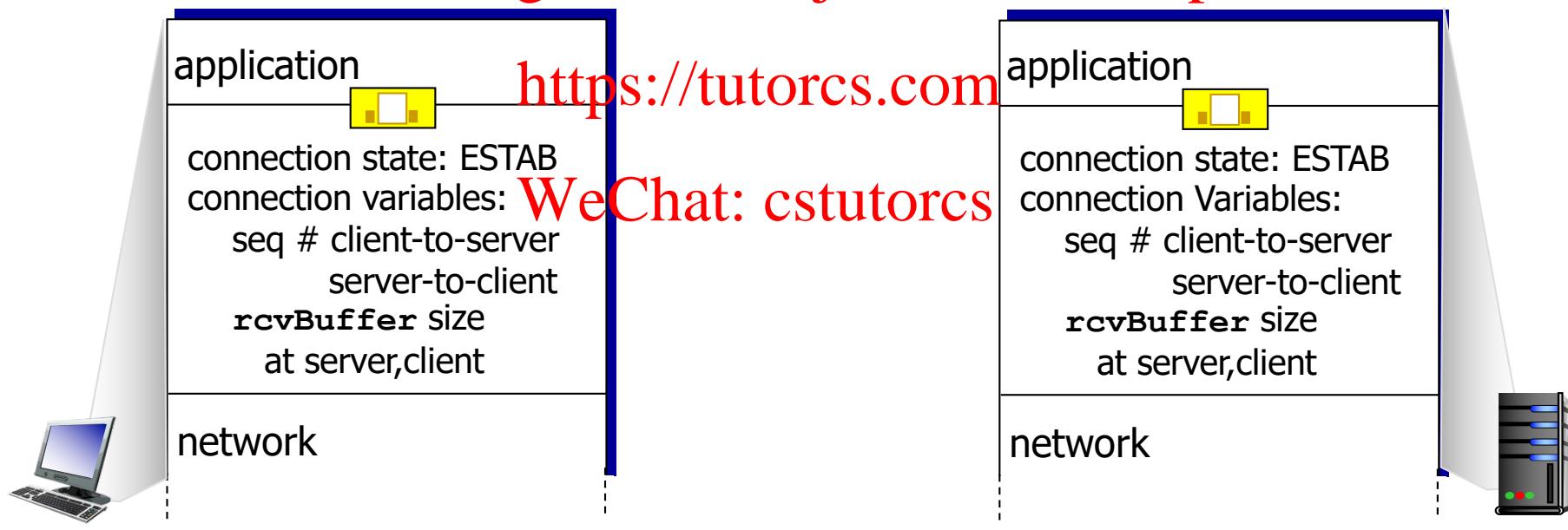
<https://tutorcs.com>
WeChat: cstutorcs

Connection Management

before exchanging data, sender/receiver “handshake”:

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

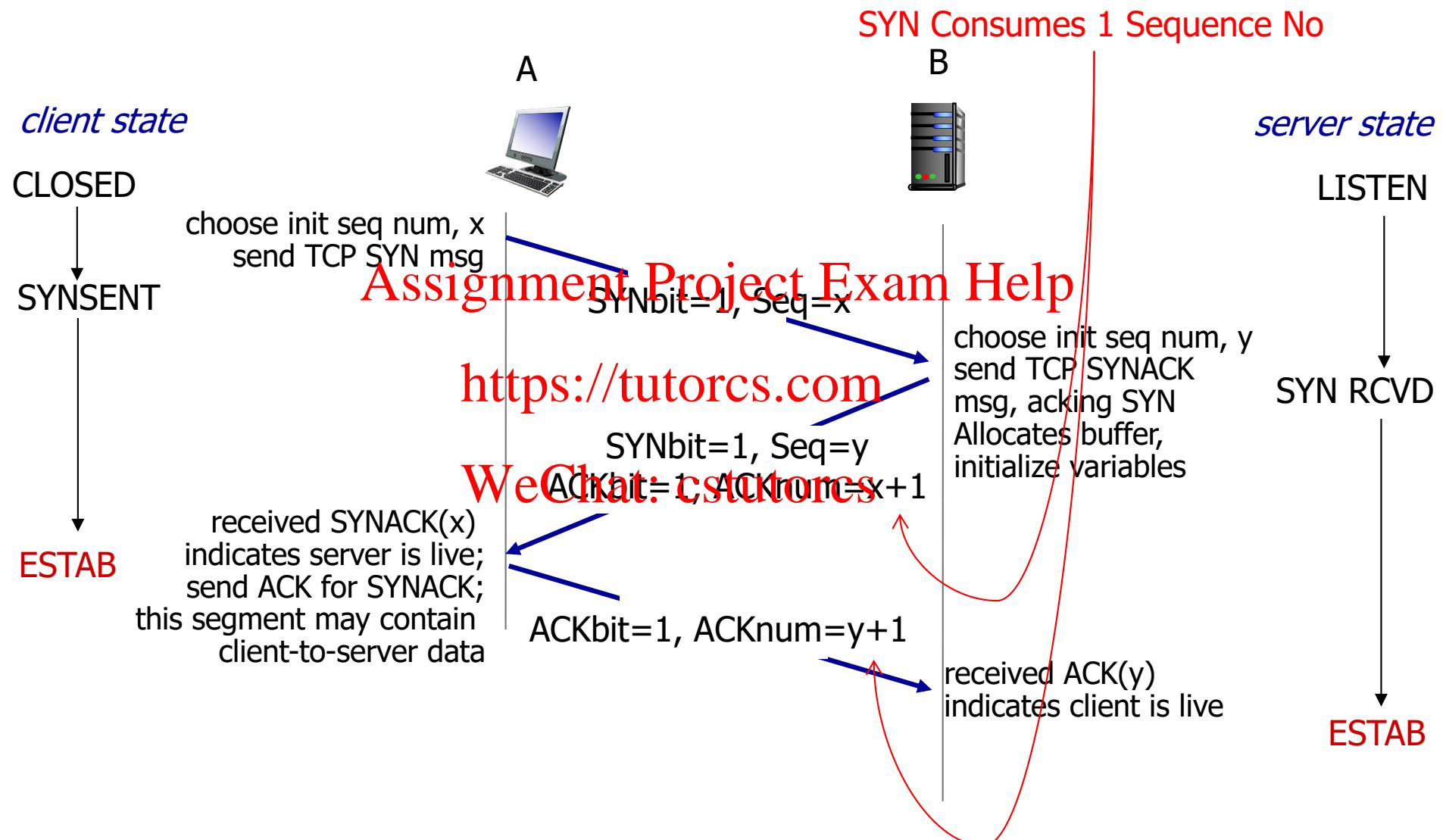
Assignment Project Exam Help



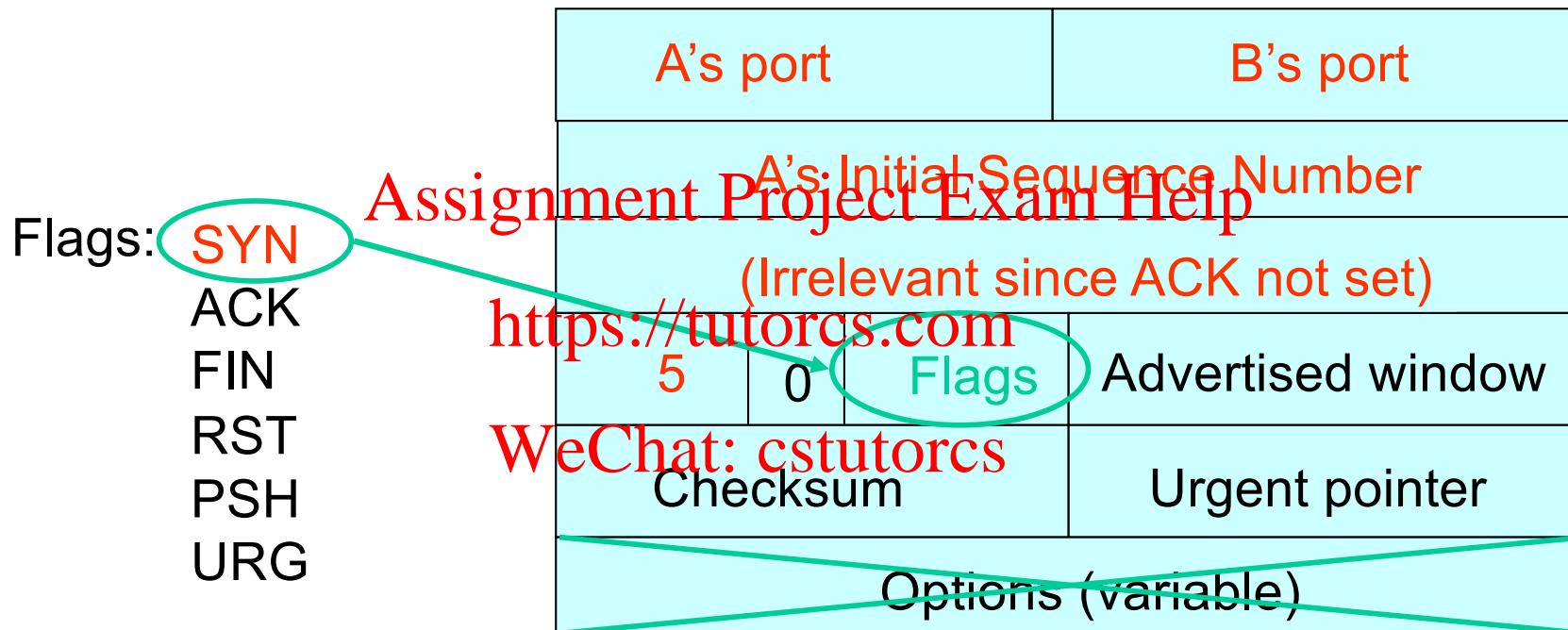
```
Socket clientSocket =  
    newSocket("hostname", "port  
    number");
```

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

TCP 3-way handshake

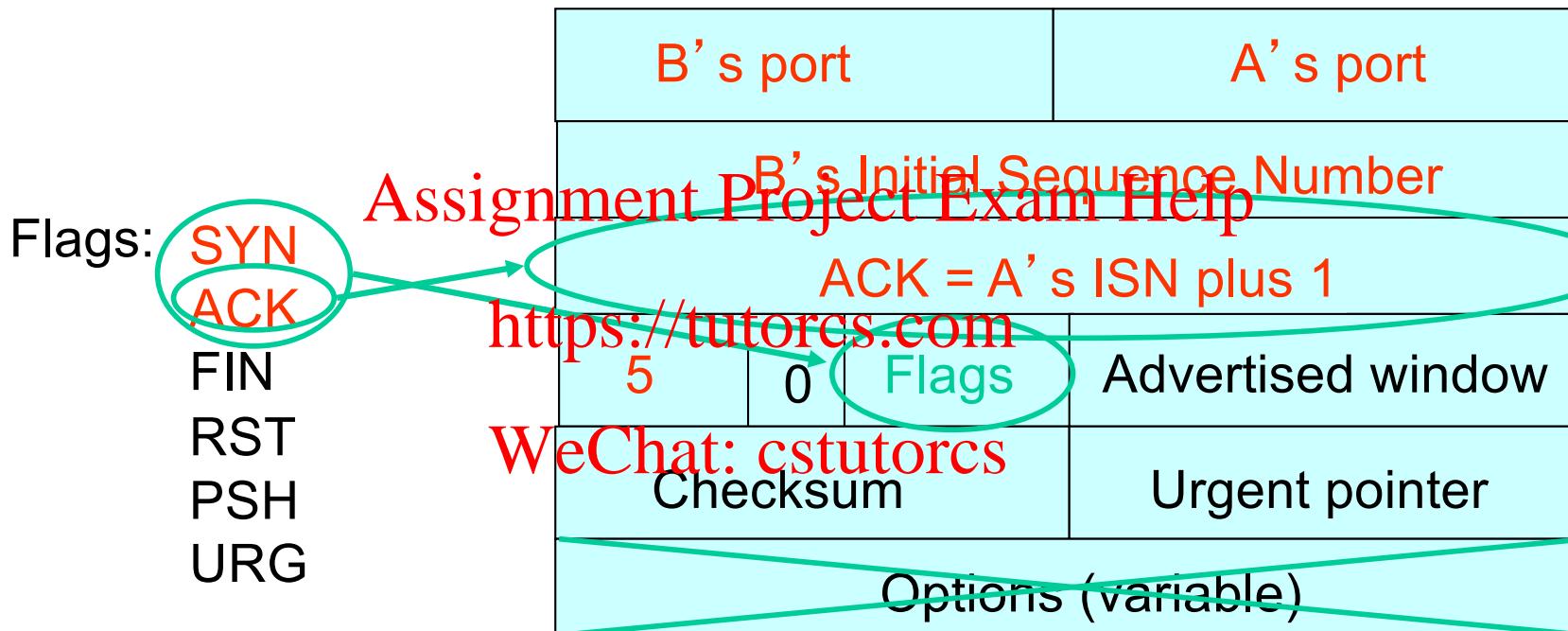


Step 1: A's Initial SYN Packet



A tells B it wants to open a connection...

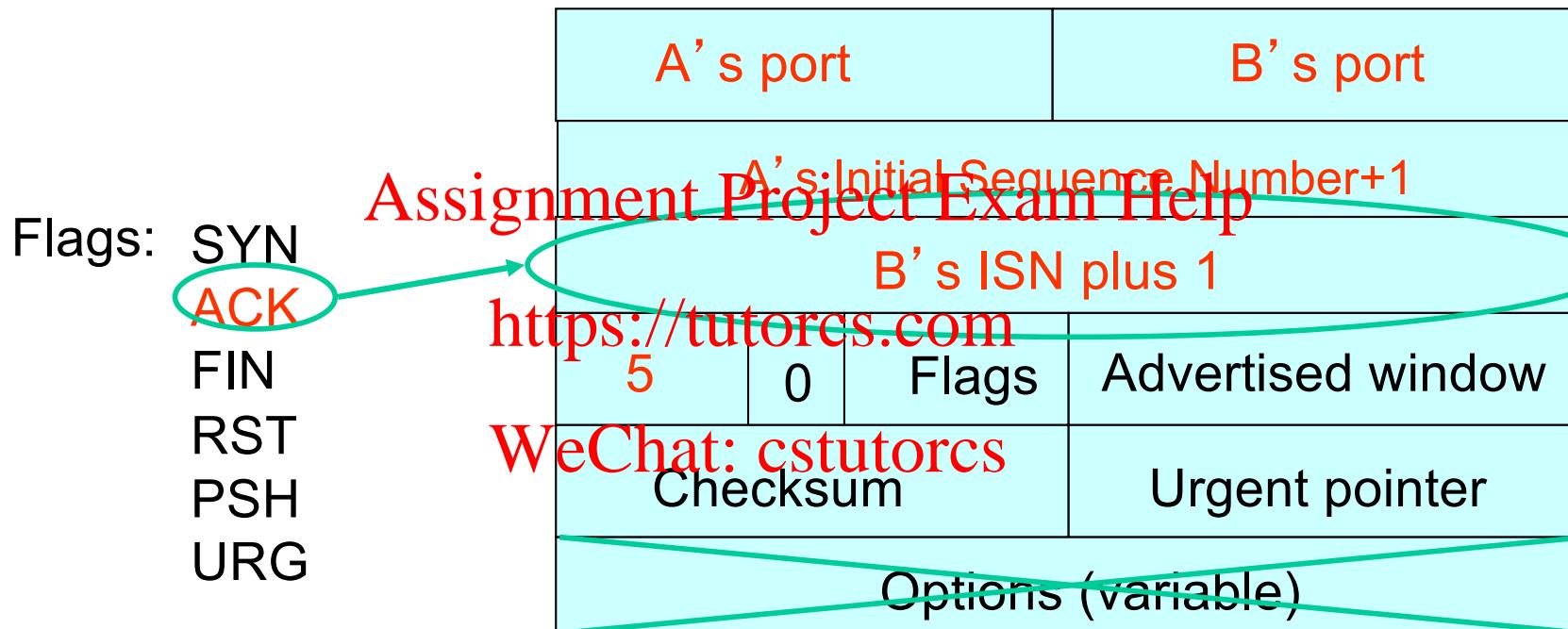
Step 2: B's SYN-ACK Packet



B tells A it accepts, and is ready to hear the next byte...

... upon receiving this packet, A can start sending data

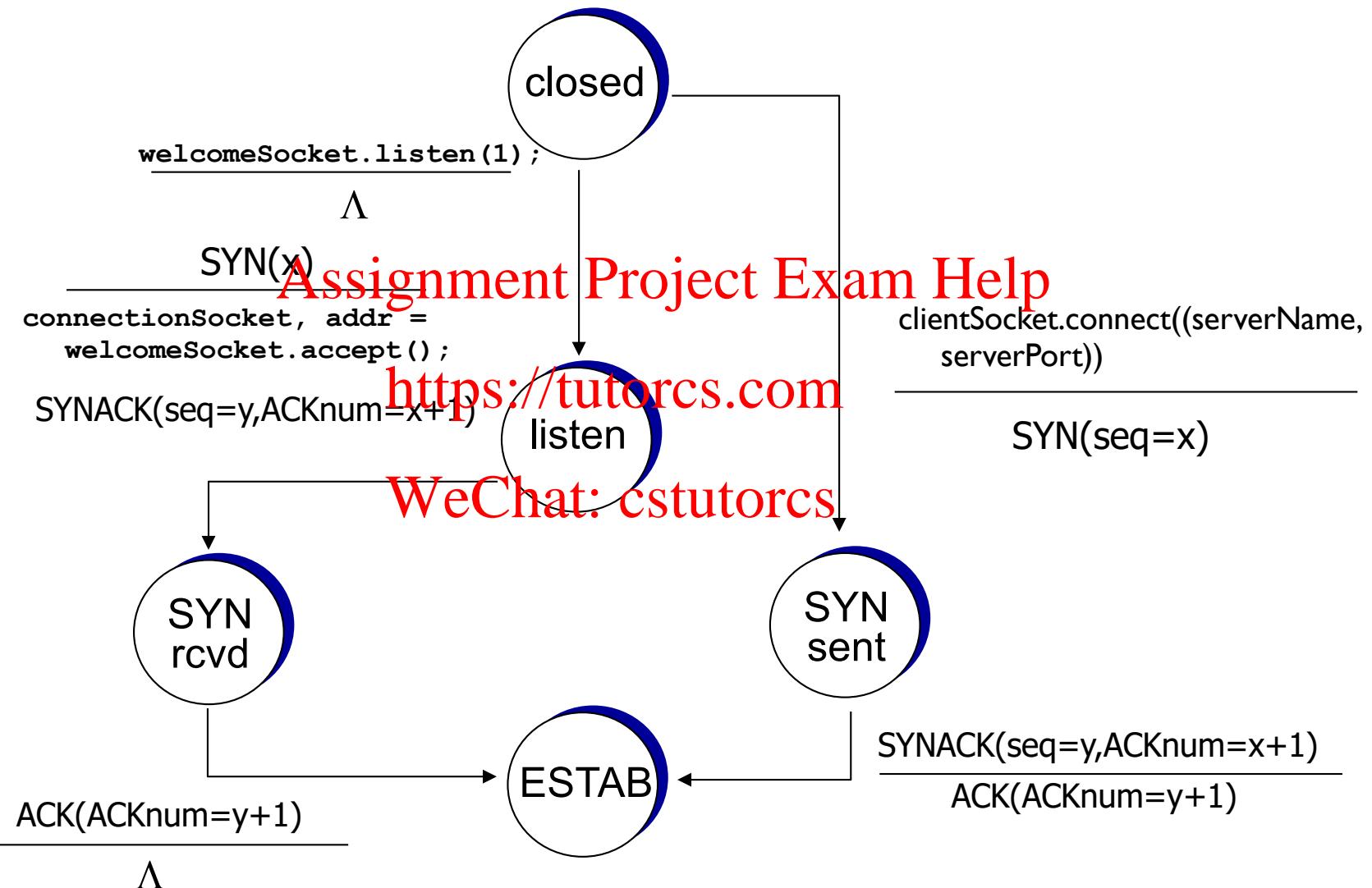
Step 3: A's ACK of the SYN-ACK



A tells B it's likewise okay to start sending

... upon receiving this packet, B can start sending data

TCP 3-way handshake: FSM



What if the SYN Packet Gets Lost?

- ❖ Suppose the SYN packet gets lost
 - Packet is lost inside the network, or:
 - Server **discards** the packet (e.g., it's too busy)
Assignment Project Exam Help
- ❖ Eventually, no SYN-ACK arrives
 - Sender sets a **timer** and **waits** for the SYN-ACK
 - ... and retransmits the SYN if needed
WeChat: cstutorcs
- ❖ How should the TCP sender set the timer?
 - Sender has **no idea** how far away the receiver is
 - Hard to guess a reasonable length of time to wait
 - **SHOULD** (RFCs 1122,2988) use default of **3 second**,
RFC 6298 use default of **1 second**

SYN Loss and Web Downloads

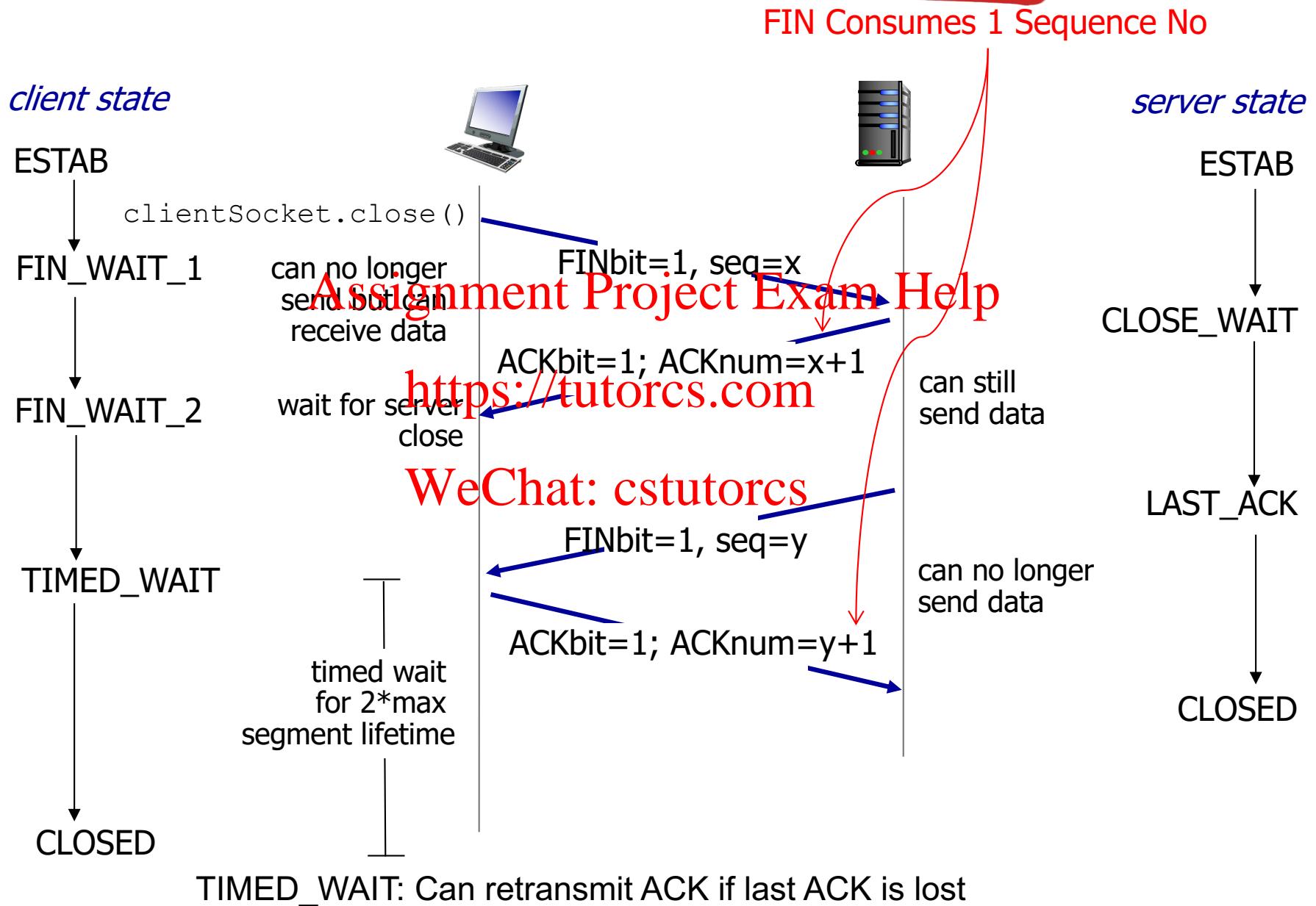
- ❖ User clicks on a hypertext link
 - Browser creates a socket and does a “connect”
 - The “connect” triggers the OS to transmit a SYN
- ❖ If the SYN ~~Assignment Project Exam Help~~
 - 1-3 seconds of delay: can be ~~very long~~ <https://tutorcs.com>
 - User may become impatient
 - ... and click the ~~WeChat; cstutorcs~~ hyperlink again, or click “reload”
- ❖ User triggers an “abort” of the “connect”
 - Browser creates a **new** socket and another “connect”
 - Essentially, forces a faster send of a new SYN packet!
 - Sometimes very effective, and the page comes quickly

TCP: closing a connection

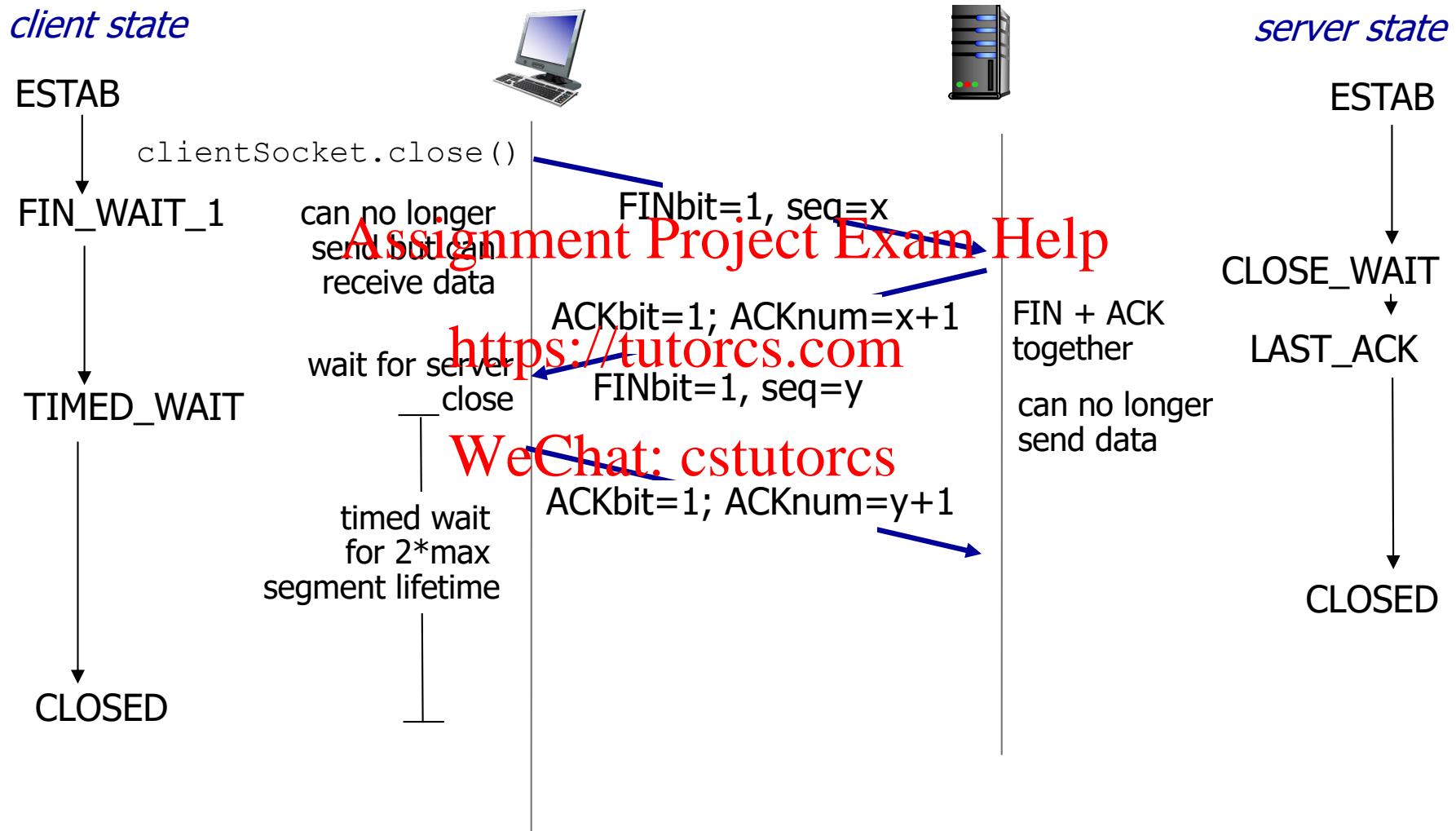
- ❖ client, server each close their side of connection
 - send TCP segment with FIN bit = 1
- ❖ respond to received FIN with ACK
 - on receiving FIN, ACK can be combined with own FIN
- ❖ simultaneous FIN exchanges can be handled
<https://tutorcs.com>

WeChat: cstutorcs

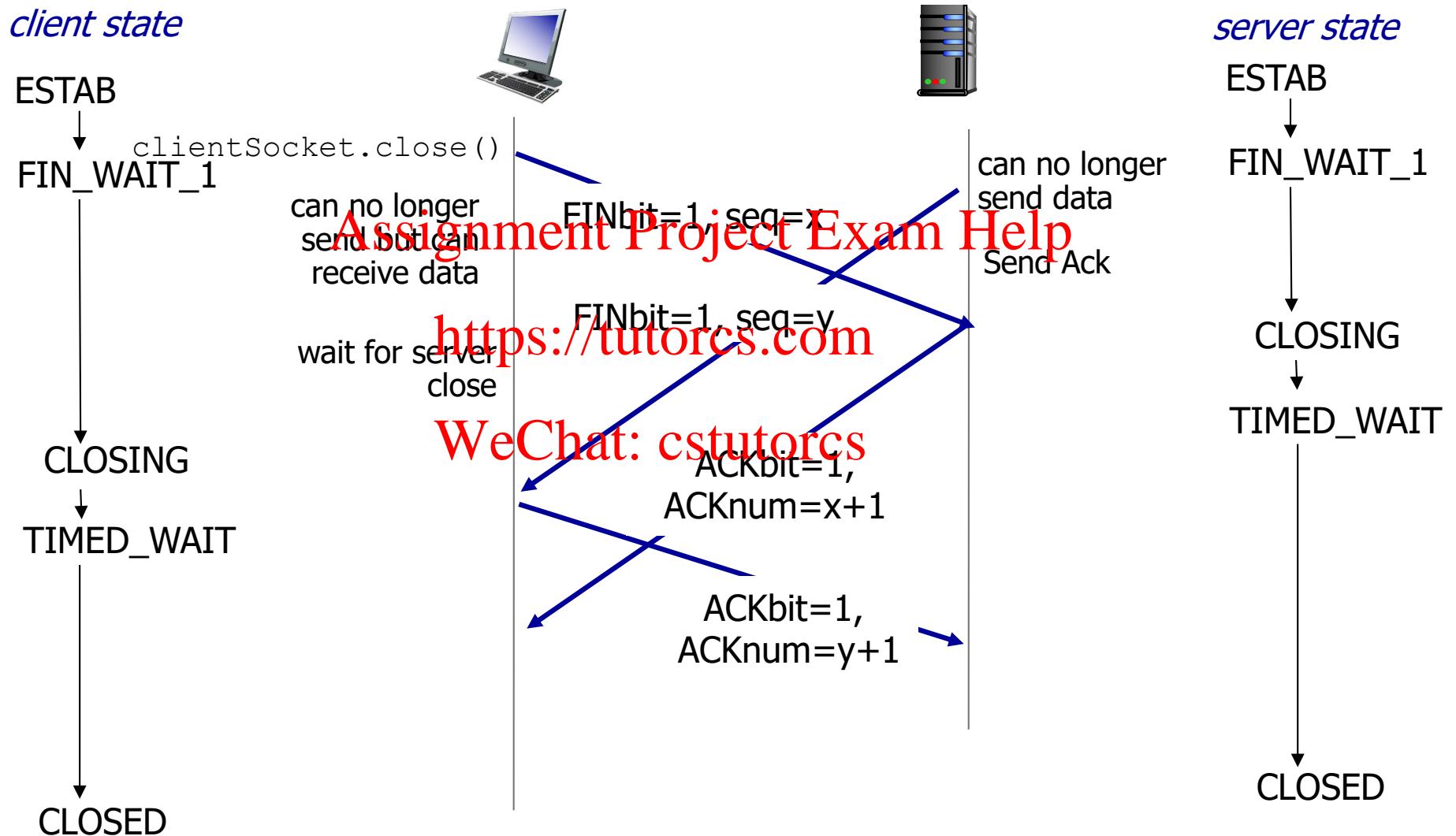
Normal Termination, One at a Time



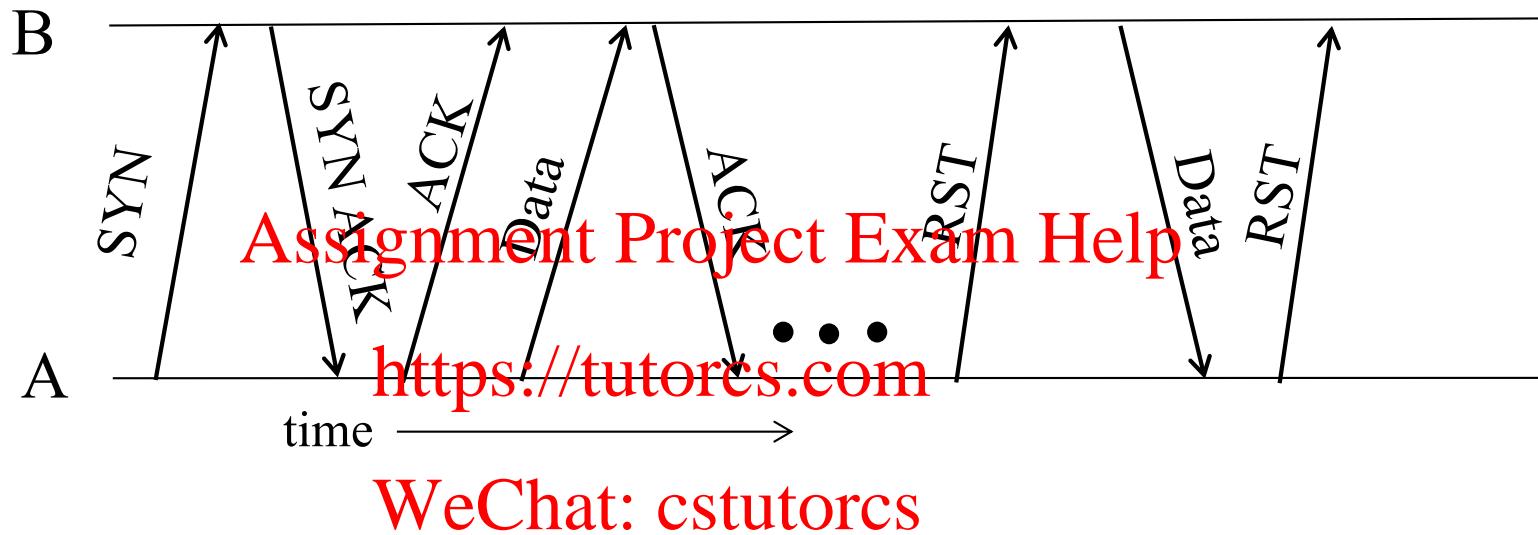
Normal Termination, Both Together



Simultaneous Closure

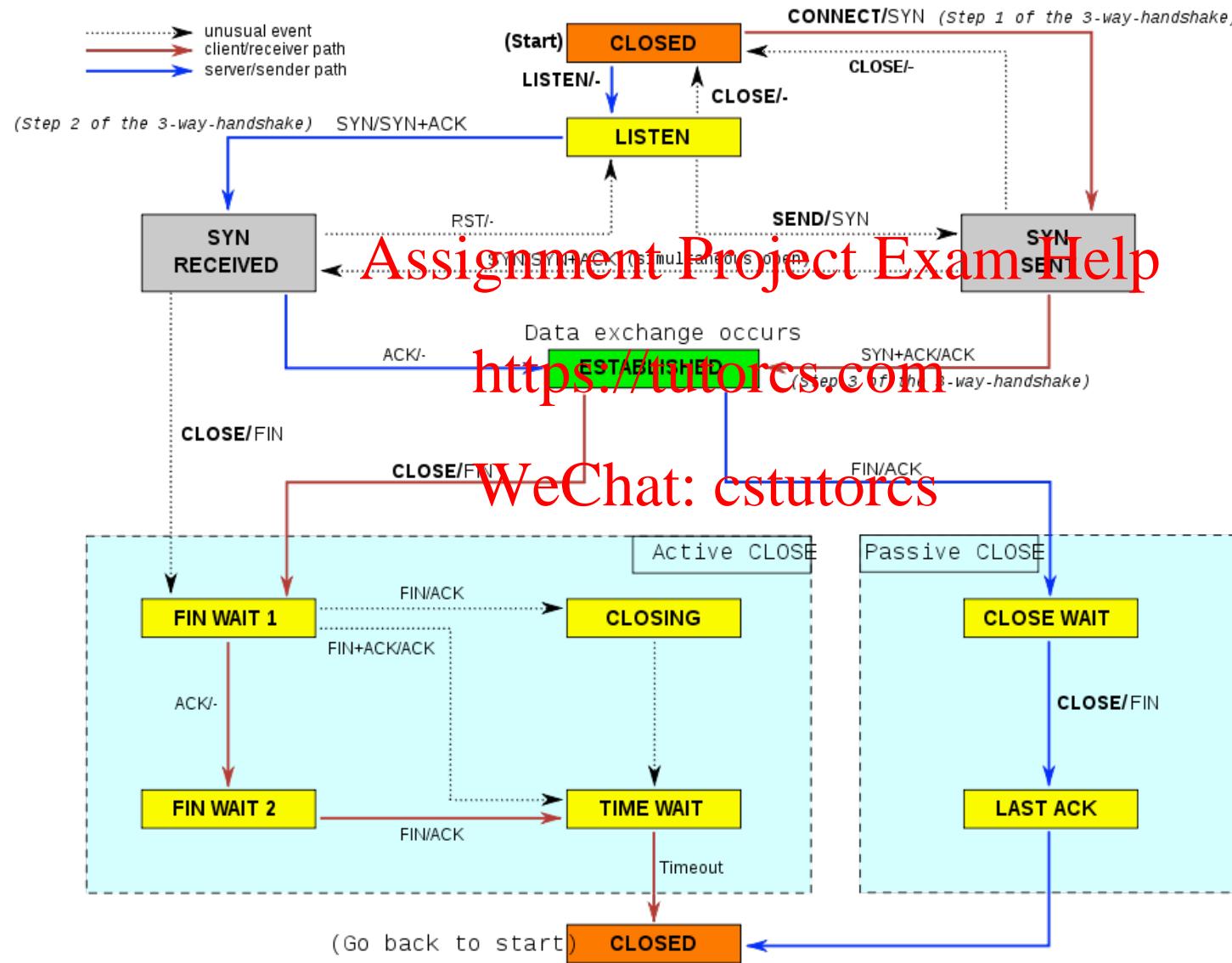


Abrupt Termination



- ❖ A sends a RESET (**RST**) to B
 - E.g., because application process on A **crashed**
- ❖ **That's it**
 - B does **not** ack the **RST**
 - Thus, **RST** is **not** delivered **reliably**
 - And: any data in flight is **lost**
 - But: if B sends anything more, will elicit **another RST**

TCP Finite State Machine



TCP SYN Attack (SYN flooding)

- ❖ Miscreant creates a fake SYN packet
 - Destination is IP address of victim host (usually some server)
 - Source is some spoofed IP address
- ❖ Victim host on receiving creates a TCP connection state i.e allocates buffers, creates variables, etc and sends SYN ACK to the spoofed address (half-open connection)
- ❖ ACK never comes back <https://tutorcs.com>
- ❖ After a timeout connection state is freed
- ❖ However for this duration the connection state is unnecessarily created
- ❖ Further miscreant sends large number of fake SYNs
 - Can easily overwhelm the victim
- ❖ Solutions:
 - Increase size of connection queue
 - Decrease timeout wait for the 3-way handshake
 - Firewalls: list of known bad source IP addresses
 - TCP SYN Cookies (explained on next slide)

TCP SYN Cookie

- ❖ On receipt of SYN, server does not create connection state
- ❖ It creates an initial sequence number (*init_seq*) that is a hash of source & dest IP address and port number of SYN packet (secret key used for hash)
 - Replies back with SYN ACK containing *init_seq*
 - Server does not need to store this sequence number
- ❖ If original SYN is genuine, an ACK will come back
 - Same hash function run on the same header fields to get the initial sequence number (*init_seq*)
 - Checks if the ACK is equal to (*init_seq+1*)
 - Only create connection state if above is true
- ❖ If fake SYN, no harm done since no state was created

<http://etherealmind.com/tcp-syn-cookies-ddos-defence/>

Quiz: TCP Connection Management?



Roughly how much time does it take for both the TCP Sender and Receiver to establish connection state since the `connect()` call?

Assignment Project Exam Help

- A. RTT
- B. 1.5RTT
- C. 2RTT
- D. 3RTT

<https://tutorcs.com>

WeChat: cstutorcs

ANSWER: B

Note that the final ACK will be typically piggybacked with the first data segment, so often the TCP connection setup is approximated to be 1RTT

Quiz: TCP Connection Management?



Assume that one end point of the TCP connection sends a FIN segment. If it never receives an ACK, what should it do?

Assignment Project Exam Help

- A. Assume that the connection is closed and do nothing <https://tutorcs.com>
- B. Retransmit the FIN
- C. Transmit an ACK
- D. Start crying

ANSWER: B

Transport Layer: 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

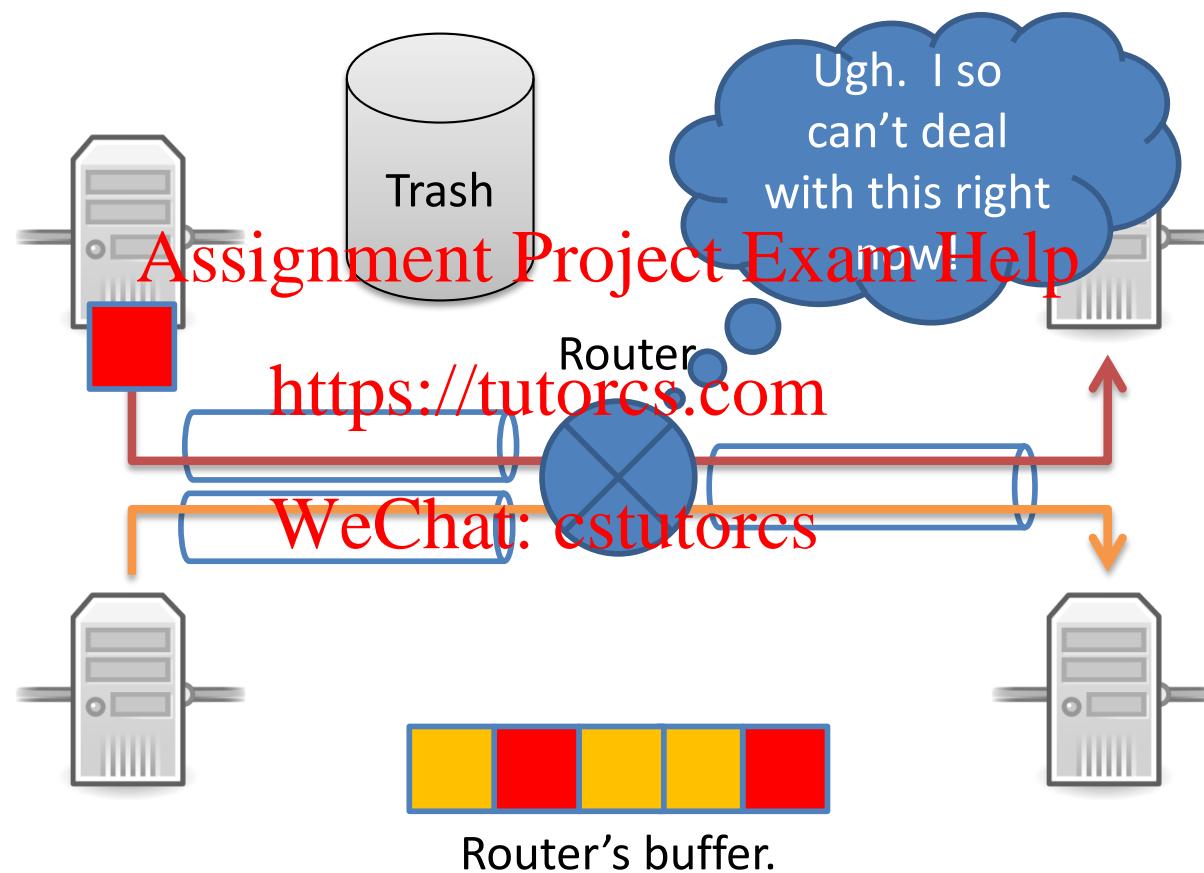
3.7 TCP congestion control

Principles of congestion control

congestion:

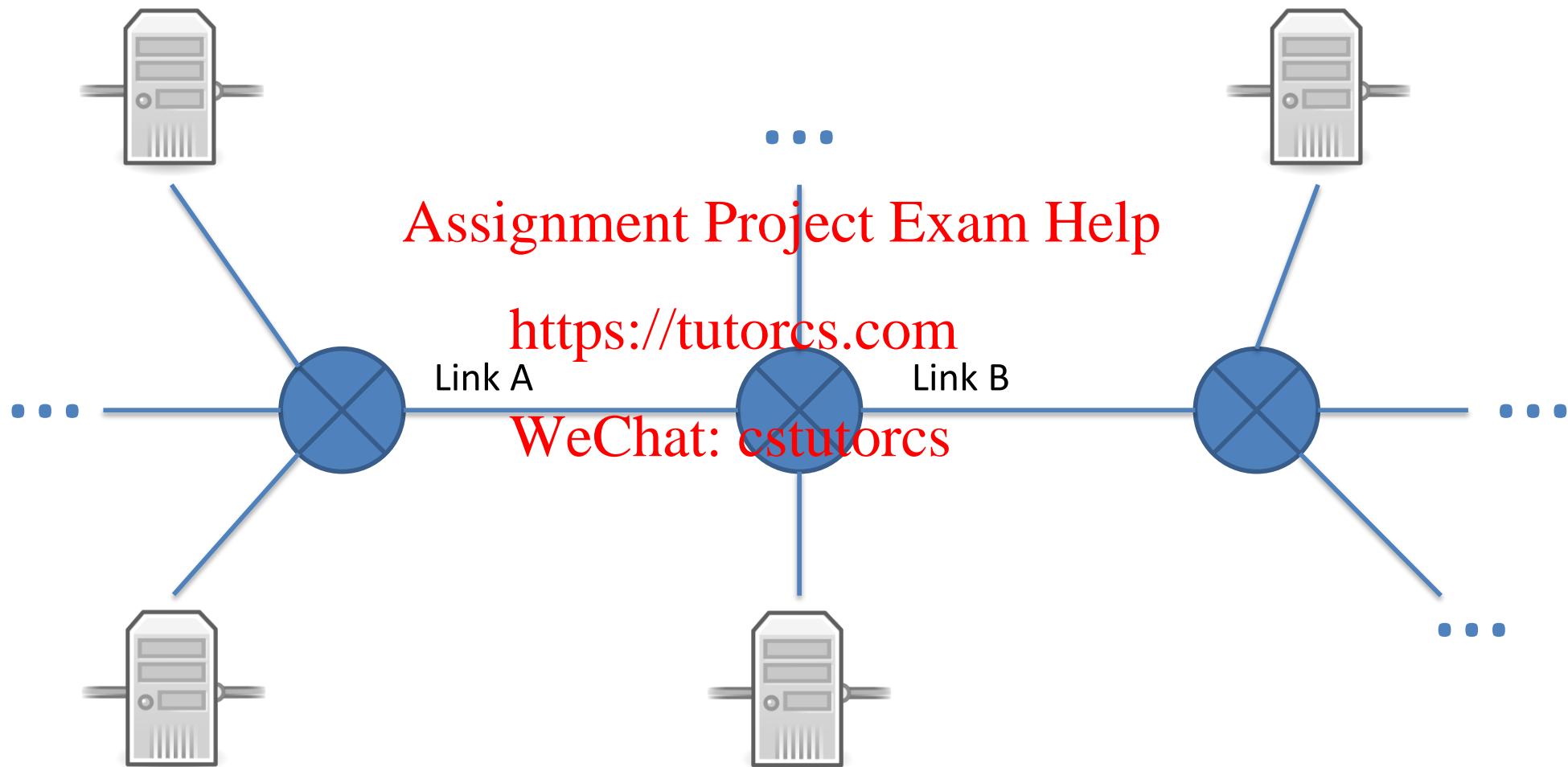
- ❖ informally: “too many sources sending too much data too fast for network to handle”
Assignment Project Exam Help
<https://tutorcs.com>
- ❖ different from flow control!
WeChat: cstutorcs
- ❖ manifestations:
 - lost packets (buffer overflow at routers)
 - long delays (queueing in router buffers)
- ❖ a top-10 problem!

Congestion

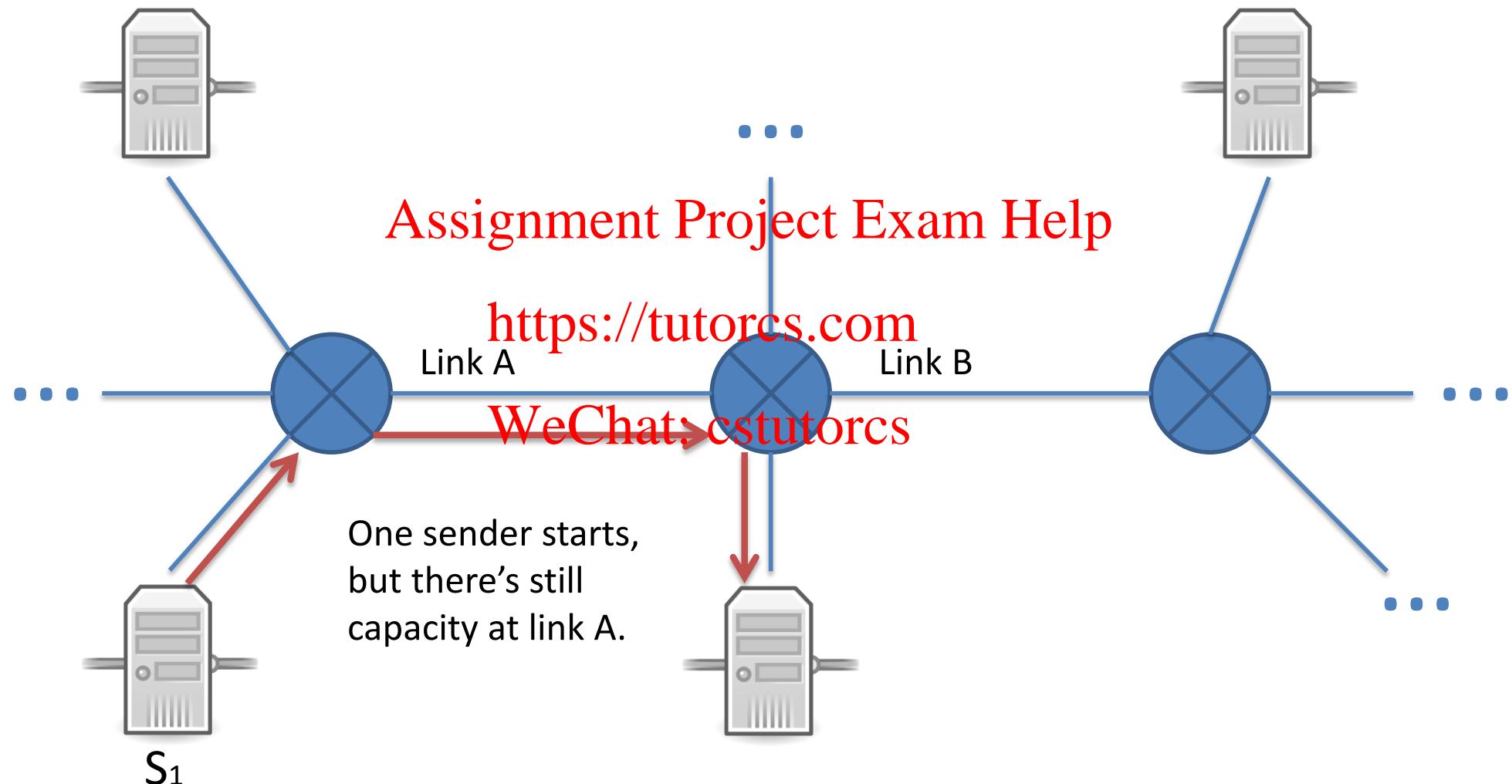


Incoming rate is faster than outgoing link can support.

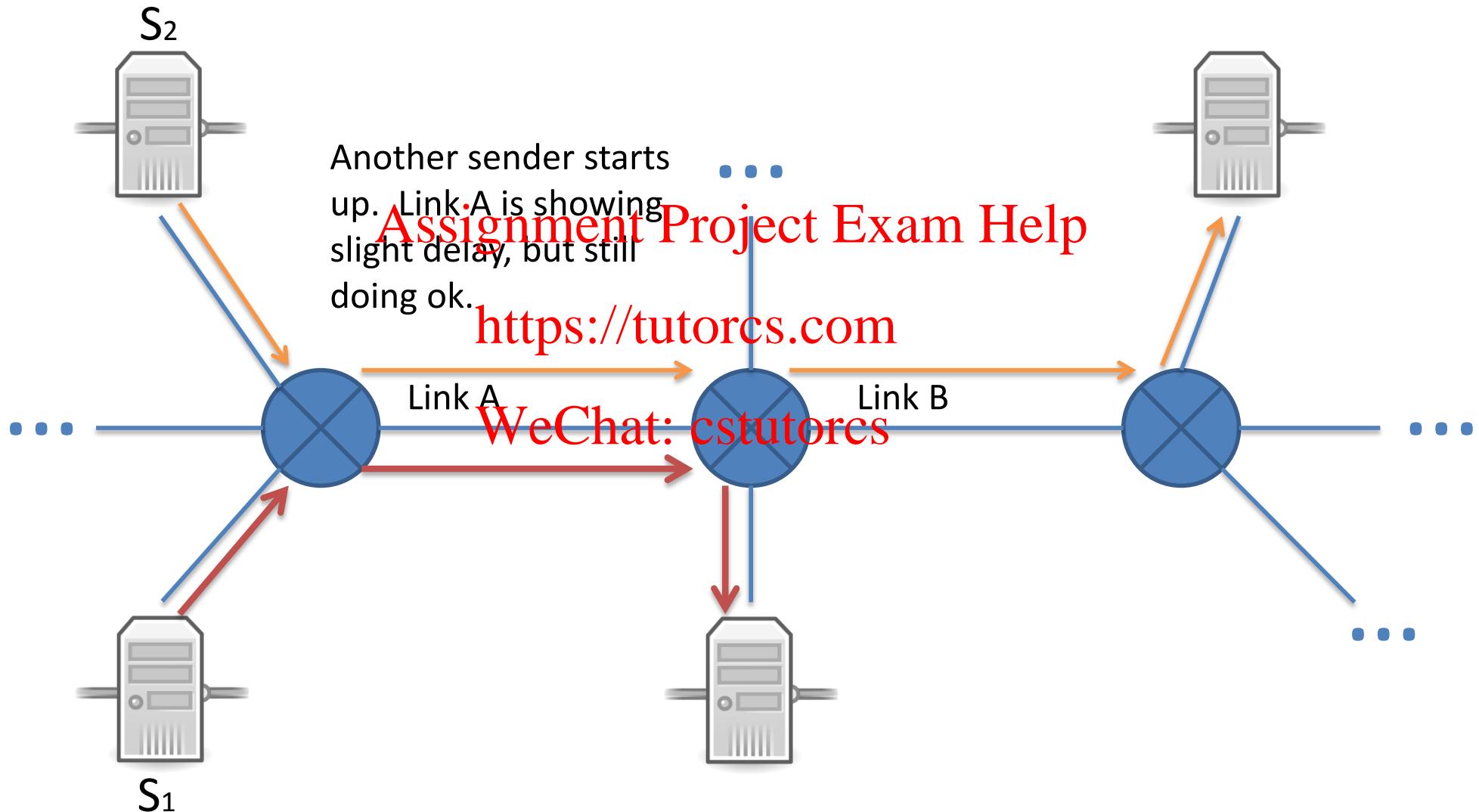
Congestion Collapse



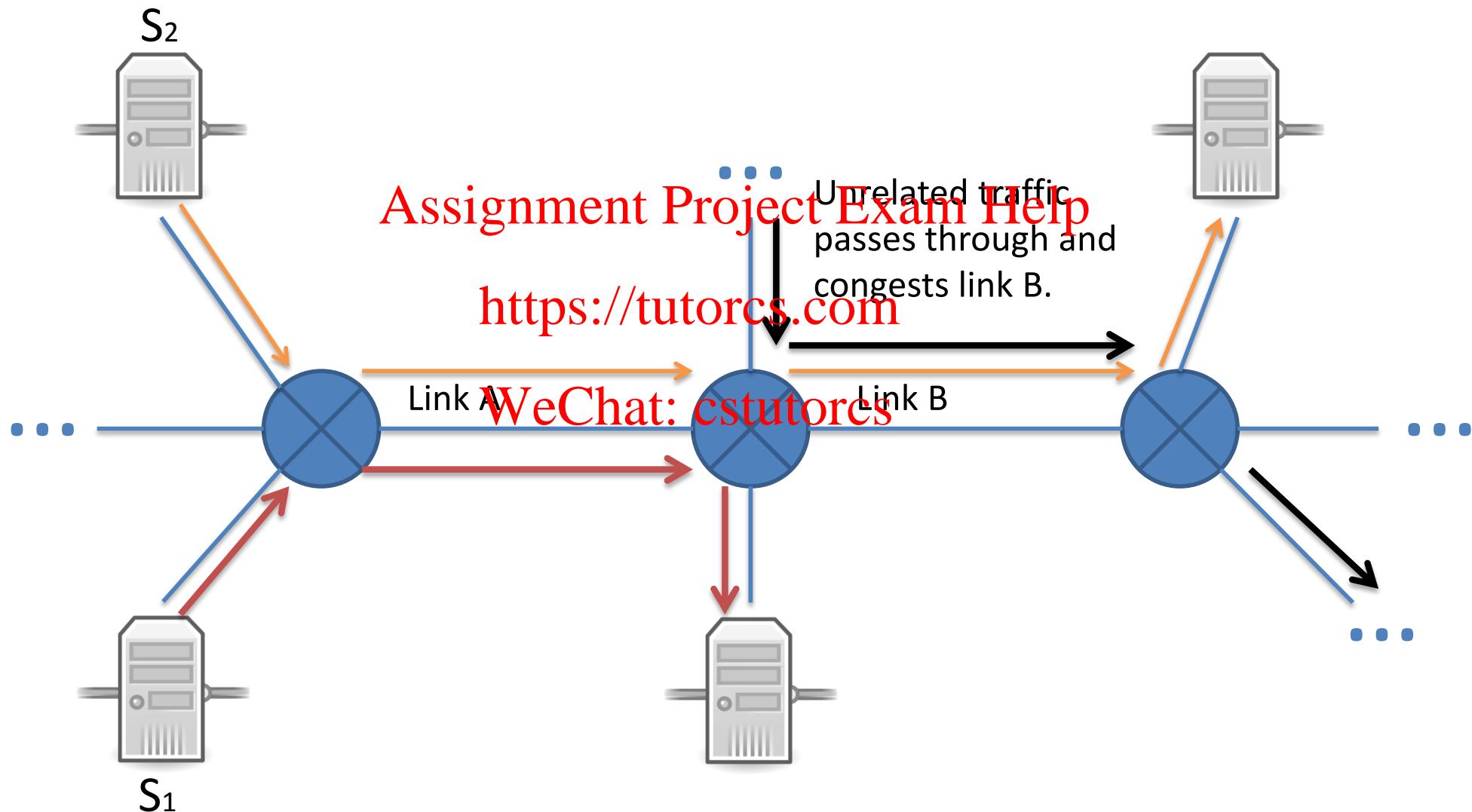
Congestion Collapse



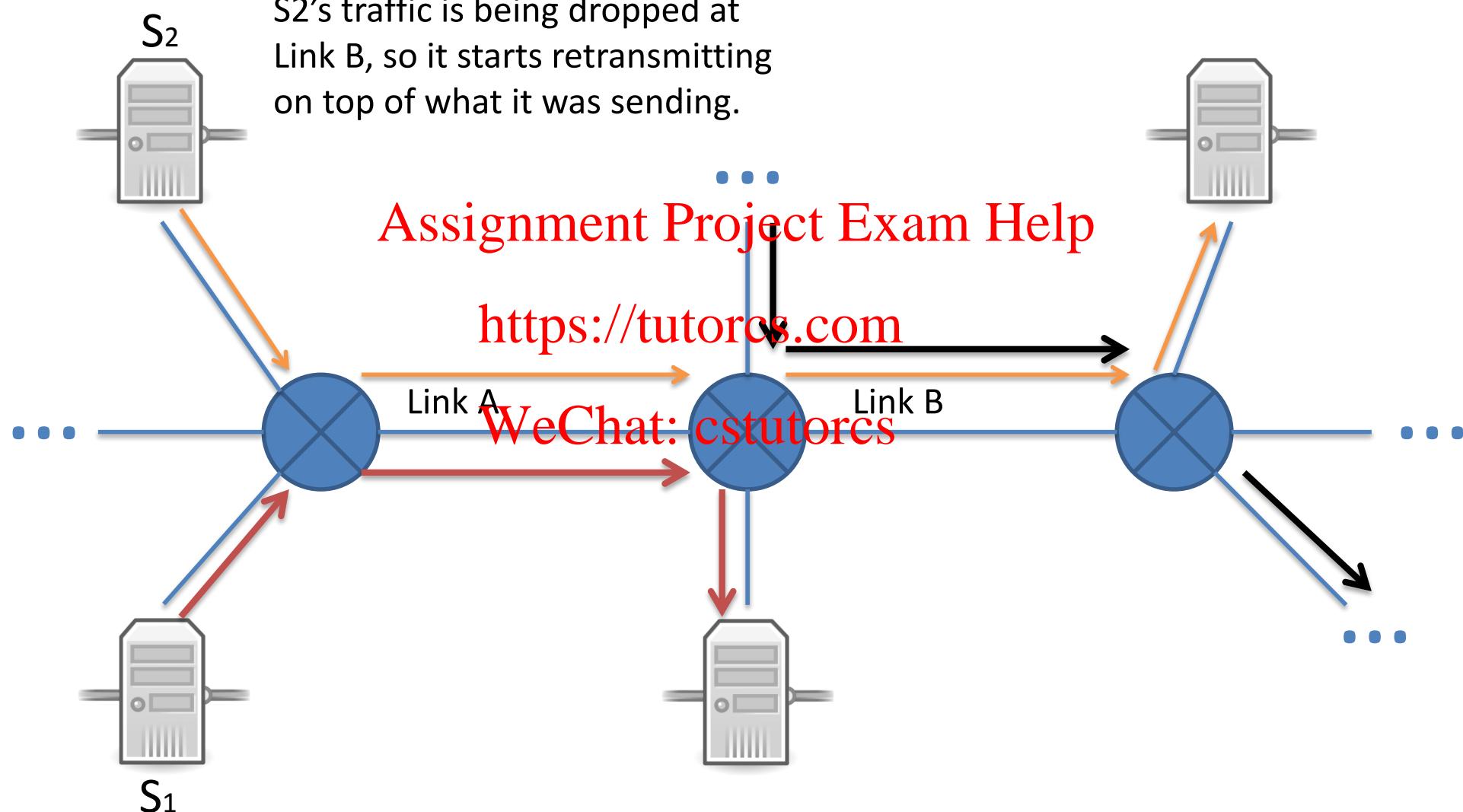
Congestion Collapse



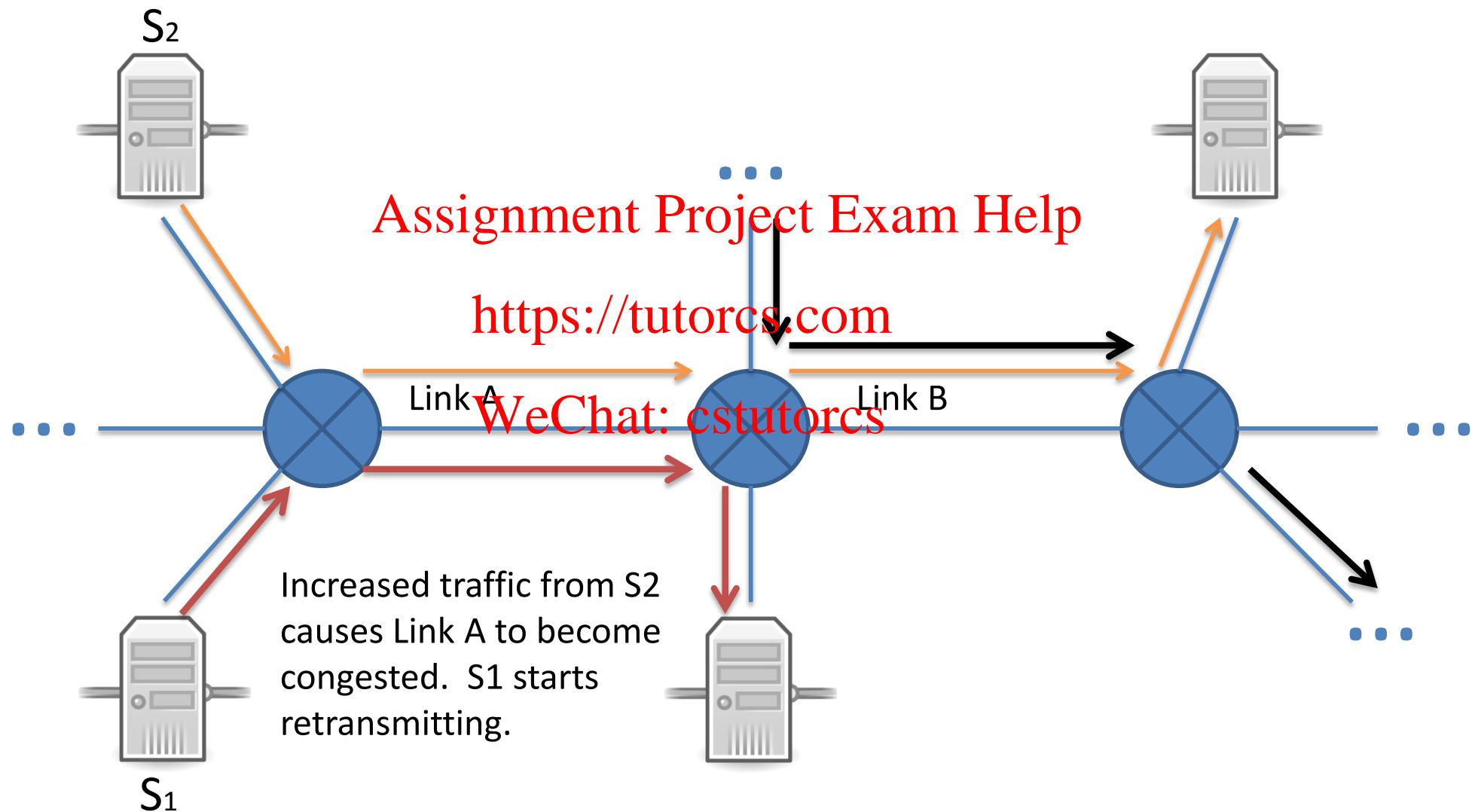
Congestion Collapse



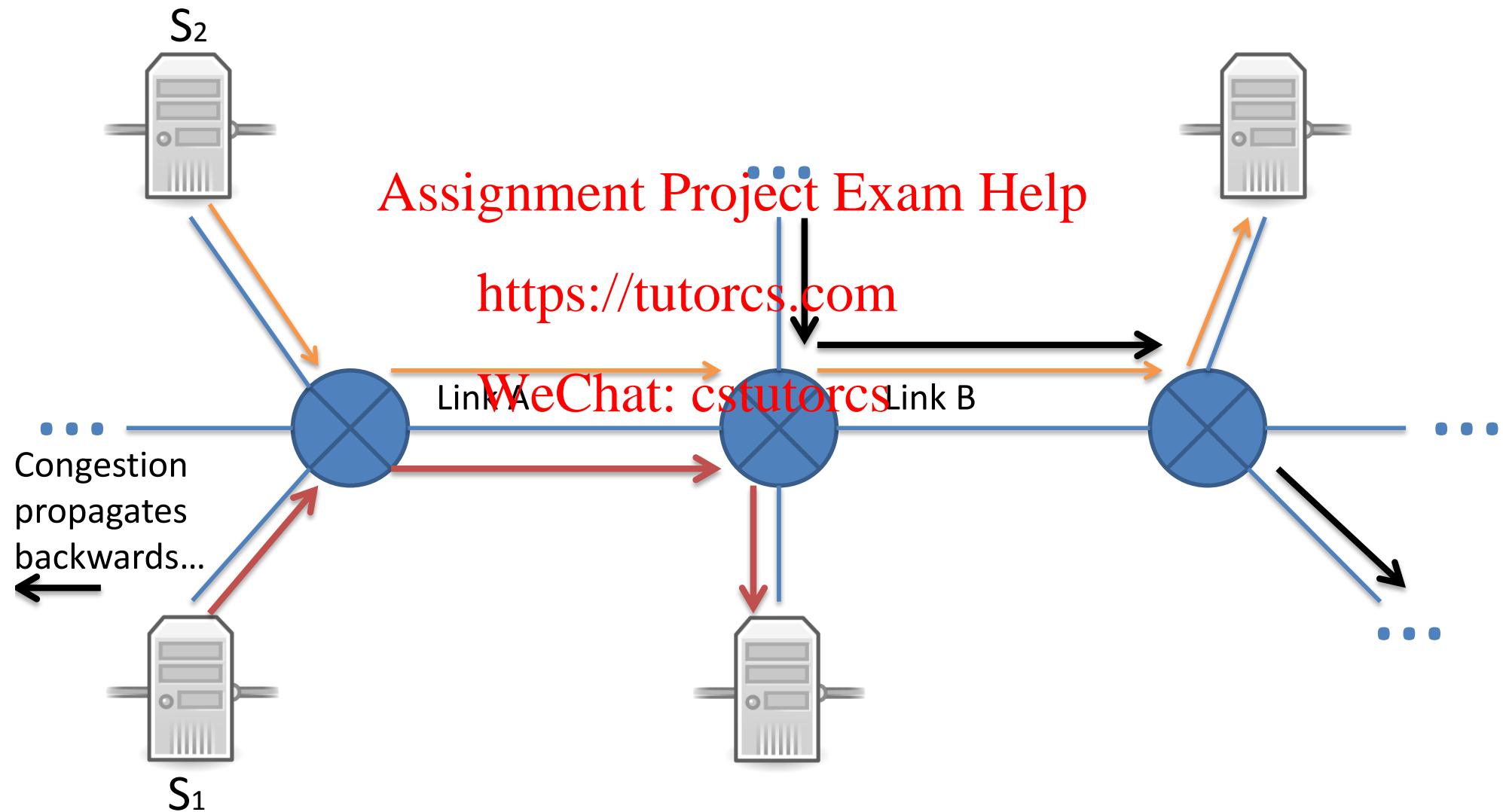
Congestion Collapse



Congestion Collapse



Congestion Collapse



Without congestion control

congestion:

- ❖ Increases delays
 - If delays > RTT, sender retransmits
- ❖ Increases loss rate
 - Dropped packets also retransmitted
- ❖ Increases retransmissions, many unnecessary
 - Wastes capacity of traffic that is never delivered
 - Increase in load results in decrease in useful work done
- ❖ Increases congestion, cycle continues ...

Cost of Congestion

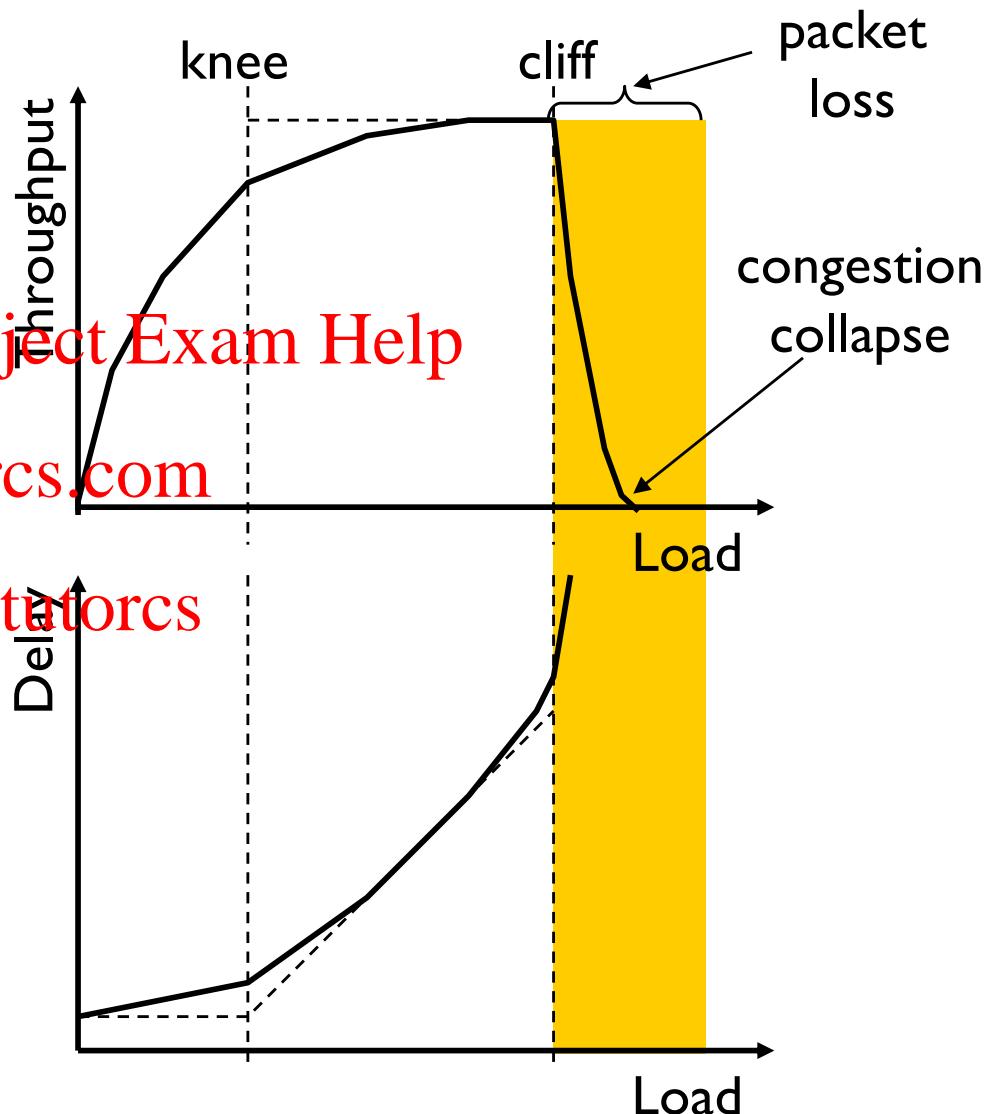
- ❖ Knee – point after which
 - Throughput increases slowly
 - Delay increases fast

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

- ❖ Cliff – point after which
 - Throughput starts to drop to zero (congestion collapse)
 - Delay approaches infinity



Congestion Collapse

This happened to the Internet (then NSFnet) in 1986

- ❖ Rate dropped from 56Kbps to 40bps
- ❖ This happened on and off for two years
- ❖ In 1988, Van Jacobson published “Congestion Avoidance and Control”
<https://tutorcs.com>
- ❖ The fix: senders voluntarily limit sending rate

Approaches towards congestion control

two broad approaches towards congestion control:

end-end congestion control:

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

network assisted congestion control:

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

Transport Layer: 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

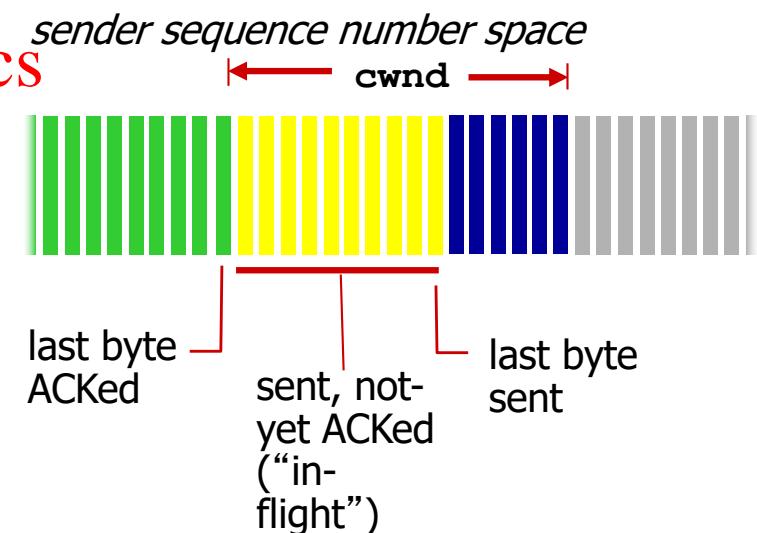
3.6 principles of congestion control

3.7 TCP congestion control

TCP's Approach in a Nutshell

- ❖ TCP connection maintains a **window**
 - Controls number of packets in flight
- ❖ *TCP sending rate*:
 - roughly: send $cwnd$ bytes, wait RTT for ACKs, then send more bytes

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



- ❖ **Vary window size to control sending rate**

All These Windows...

- ❖ Congestion Window: **CWND**
 - How many bytes can be sent without overflowing routers
 - Computed by the sender using Congestion control algorithm
- ❖ Flow control window: **Advertised / Receive Window (RWND)**
 - How many bytes can be sent without overflowing receiver's buffers
 - Determined by the receiver and reported to the sender
- ❖ Sender-side window = **minimum{CWND, RWND}**
 - Assume for this discussion that RWND >> CWND

CWND

- ❖ This lecture will talk about CWND in units of MSS
 - (Recall MSS: Maximum Segment Size, the amount of payload data in a TCP Packet)
 - This is only for pedagogical purposes
<https://tutorcs.com>

WeChat: cstutorcs

- ❖ Keep in mind that real implementations maintain CWND in bytes

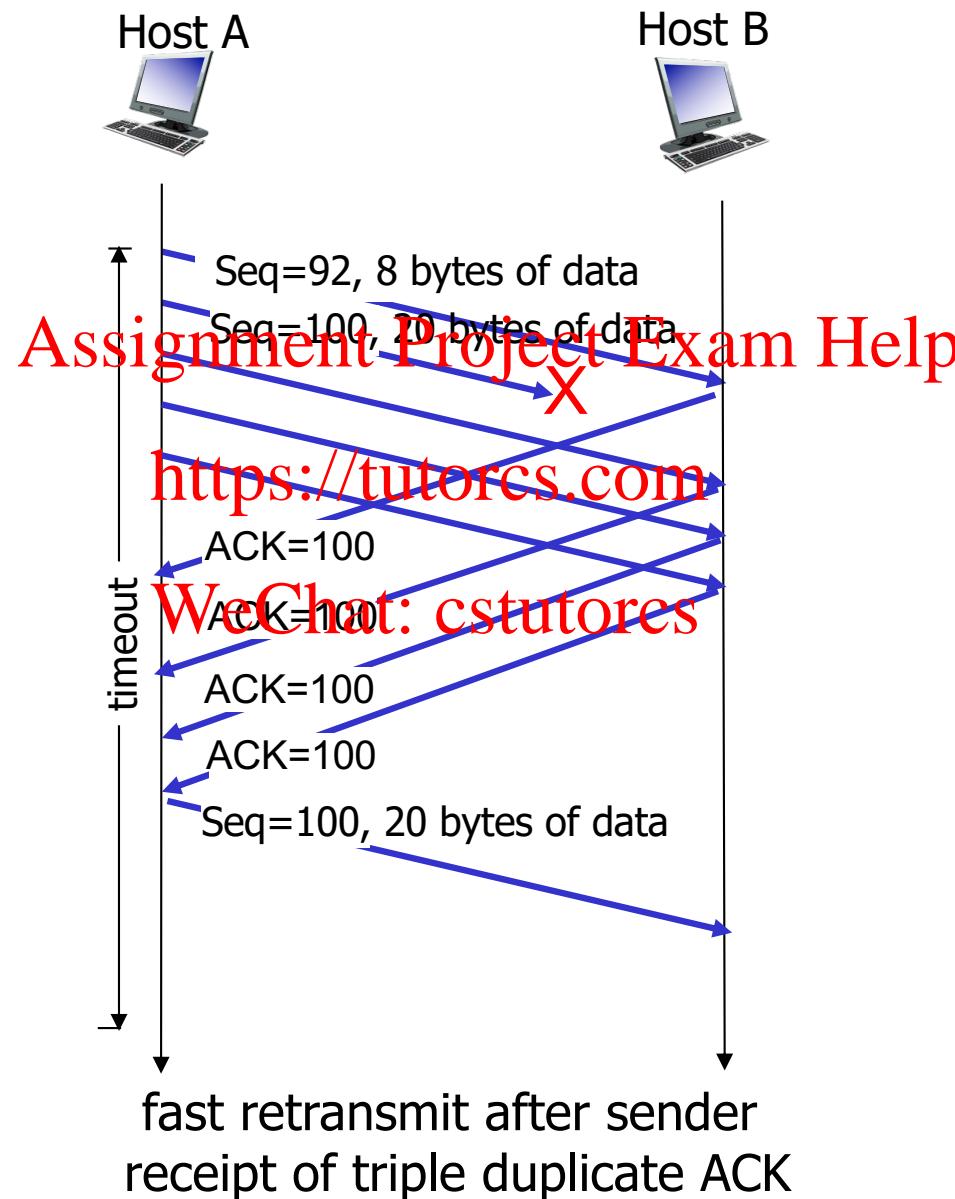
Two Basic Questions

- ❖ How does the sender detect congestion?
Assignment Project Exam Help
<https://tutorcs.com>
WeChat: cstutorcs
- ❖ How does the sender adjust its sending rate?

Detection Congestion: Infer Loss

- ❖ Duplicate ACKs: isolated loss
 - dup ACKs indicate network capable of delivering some segments
- ❖ Assignment Project Exam Help
 - <https://tutorcs.com>
- ❖ Timeout: much more serious
 - WeChat: cstutorcs
 - Not enough dup ACKs
 - Must have suffered several losses
- ❖ Will adjust rate differently for each case

RECAP: TCP fast retransmit (dup acks)



Rate Adjustment

- ❖ Basic structure:
 - Upon receipt of ACK (of new data): increase rate
Assignment Project Exam Help
<https://tutorcs.com>
 - Upon detection of loss: decrease rate
- ❖ How we increase/decrease the rate depends on the phase of congestion control we're in:
 - Discovering available bottleneck bandwidth vs.
 - Adjusting to bandwidth variations

TCP Slow Start (Bandwidth discovery)

- ❖ when connection begins, increase rate **exponentially** until **first loss event**:

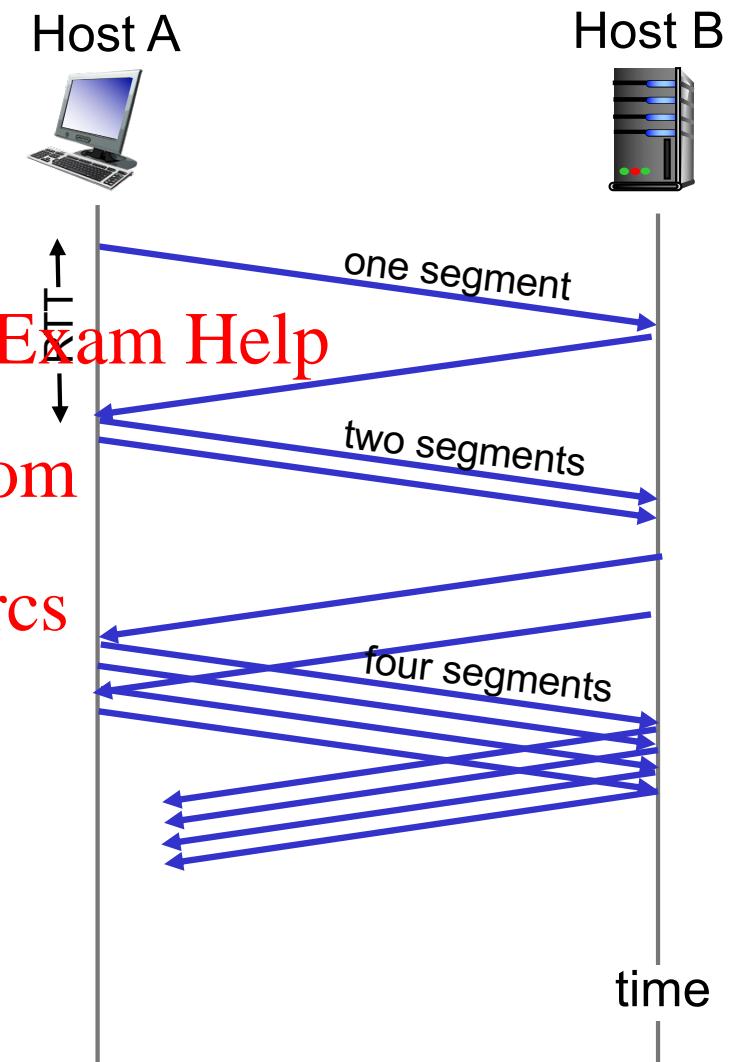
- initially $cwnd = 1 \text{ MSS}$

- double **cwnd** every RTT (full ACKs)

- Simpler implementation achieved by incrementing **cwnd** for every ACK received

- $cwnd += 1$ for each ACK

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



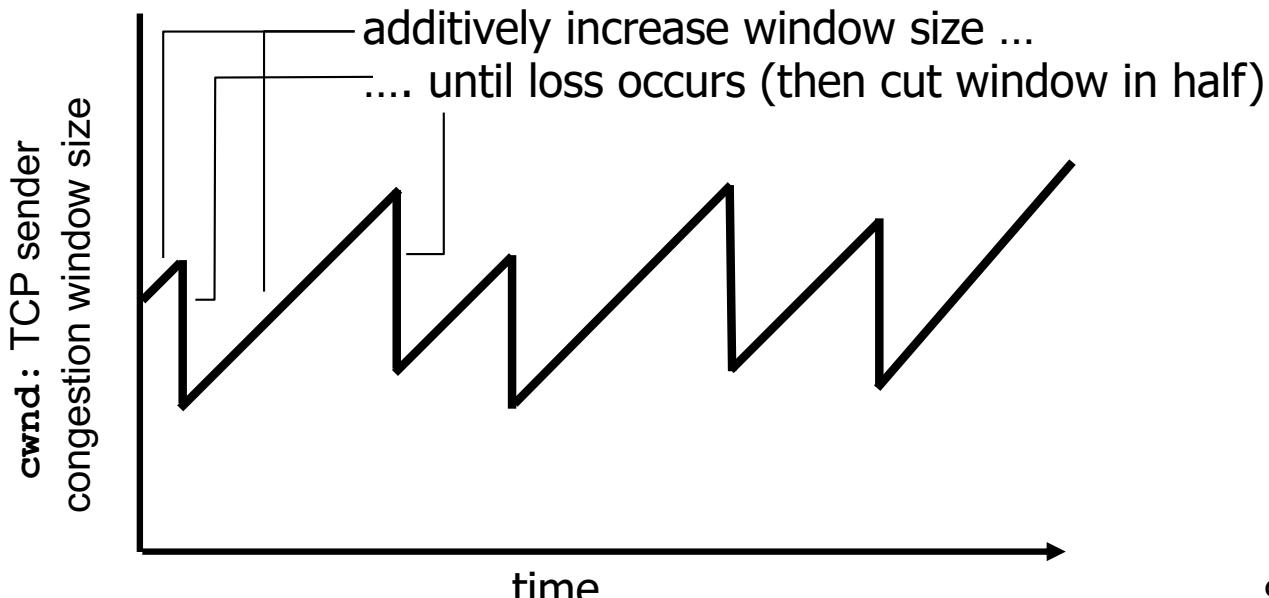
Adjusting to Varying Bandwidth

- ❖ Slow start gave an estimate of available bandwidth
- ❖ Now, want to track variations in this available bandwidth, oscillating around its current value
 - Repeated probing (rate increase) and backoff (rate decrease)
 - Known as Congestion Avoidance (CA)
- ❖ TCP uses: “Additive Increase Multiplicative Decrease” (AIMD)

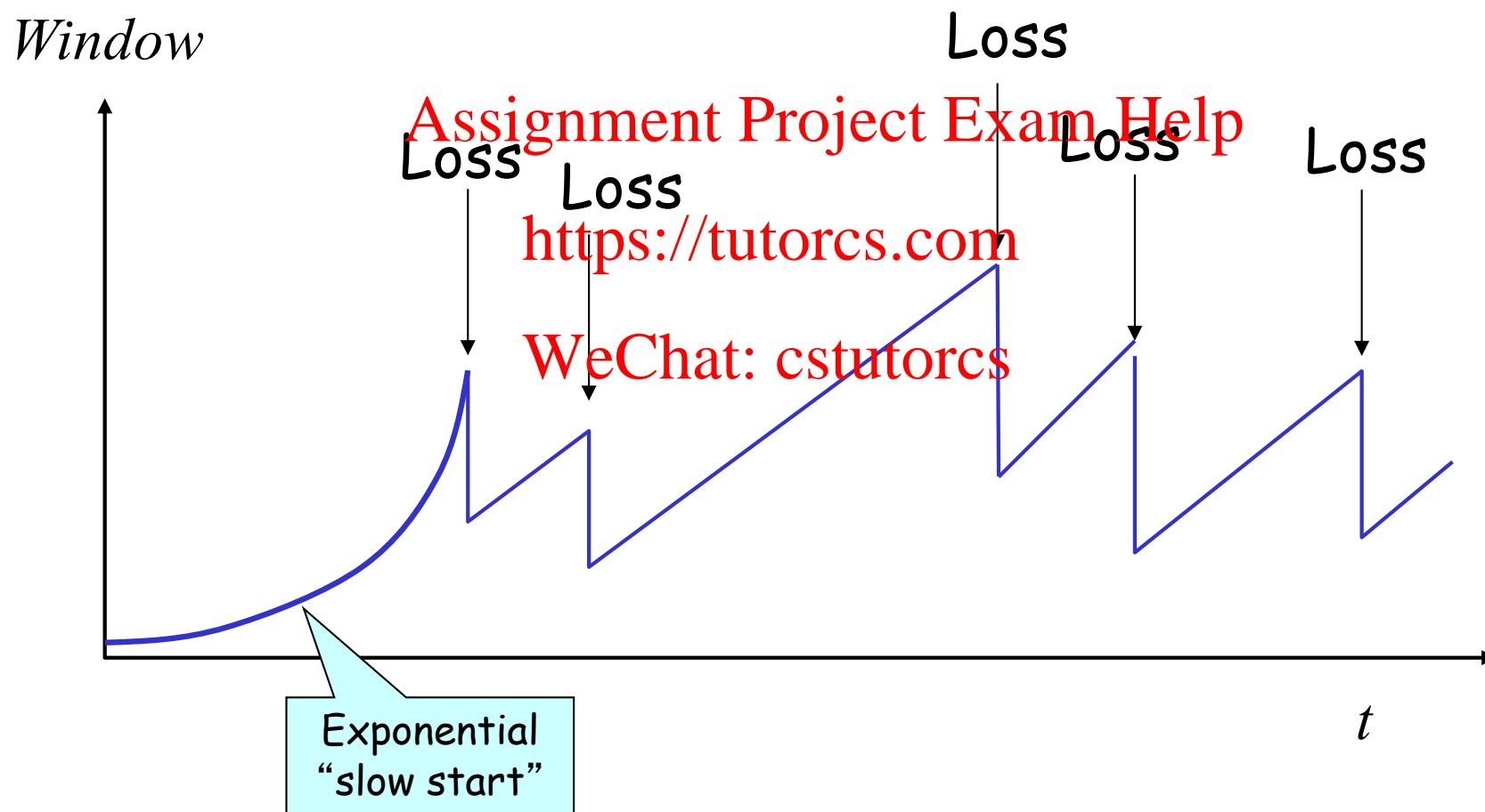
AIMD

- ❖ **approach:** sender increases transmission rate (window size), probing for usable bandwidth, until another congestion event occurs
 - **additive increase:** increase $cwnd$ by 1 MSS every RTT until loss detected
 - For each successful RTT (all ACKS), $cwnd = cwnd + 1$
 - Simple implementation: for each ACK, $cwnd = cwnd + 1/cwnd$ (since there are $cwnd/MSS$ packets in a window)
 - **multiplicative decrease:** cut $cwnd$ in half after loss

AIMD saw tooth behavior: probing for bandwidth



Leads to the TCP “Sawtooth”



Slow-Start vs. AIMD

- ❖ When does a sender stop Slow-Start and start Additive Increase?

Assignment Project Exam Help
<https://tutorcs.com>
- ❖ Introduce a “slow start threshold” (`ssthresh`)
 - Initialized to a large value

WeChat: cstutorcs
- ❖ Convert to AI when $cwnd = ssthresh$, sender switches from slow-start to AIMD-style increase
 - On timeout, $ssthresh = CWND/2$

Implementation

❖ State at sender

- CWND (initialized to a small constant)
- ssthresh (initialized to a large constant)
- [Also dupACKcount and timer, as before]
<https://tutorcs.com>

❖ Events

WeChat: cstutorcs

- ACK (new data)
- dupACK (duplicate ACK for old data)
- Timeout

Event: ACK (new data)

- ❖ If $CWND < ssthresh$

- $CWND += +$

- Hence after one RTT (All ACKs with no drops):

$$CWND = 2 \times CWND$$

<https://tutorcs.com>

WeChat: cstutorcs

Event: ACK (new data)

- ❖ If $CWND < ssthresh$
 - $CWND += 1$
- ❖ Else
 - $CWND = \min(CWND + 1, \frac{I}{CWND})$

Slow start phase

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

*“Congestion
Avoidance” phase
(additive increase)*

- Hence after one RTT (All ACKs with no drops):

$$CWND = CWND + 1$$

Event: dupACK

- ❖ dupACKcount ++
Assignment Project Exam Help
- ❖ If dupACKcount = 3 /* fast retransmit */
 - ssthresh = CWND/2
WeChat: cstutorcs
 - CWND = CWND/2

Event: TimeOut

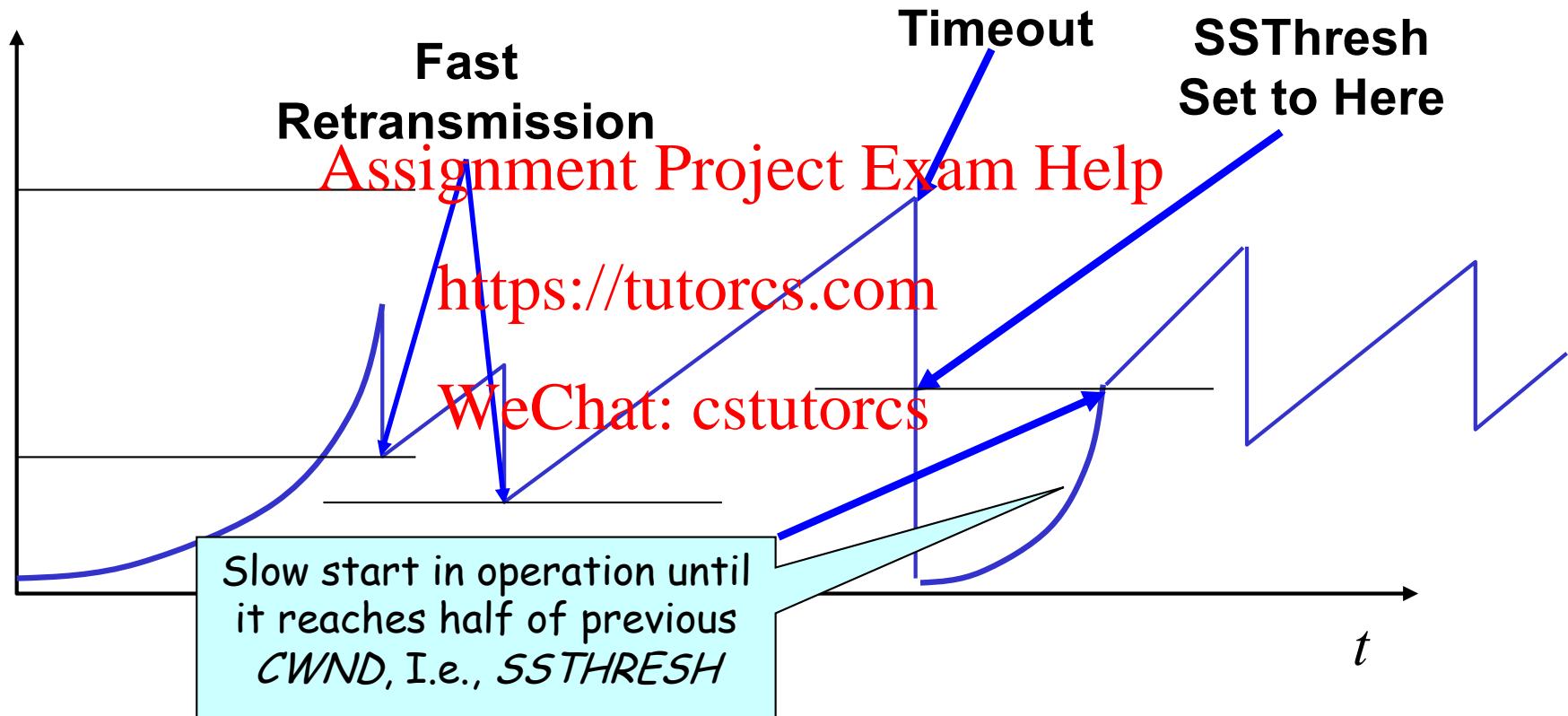
❖ On Timeout

- $ssthresh \leftarrow CWND/2$
- $CWND \leftarrow https://tutorcs.com$

WeChat: cstutorcs

Example

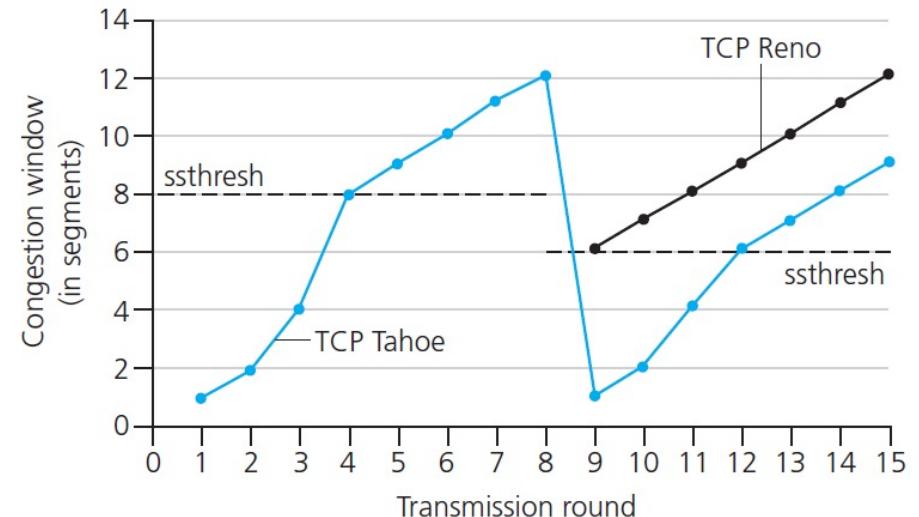
Window



Slow-start restart: Go back to $CWND = 1$ MSS, but take advantage of knowing the previous value of $CWND$

TCP Flavours

- ❖ TCP-Tahoe
 - $cwnd = 1$ on triple dup ACK & timeout
- ❖ TCP-Reno
 - $cwnd = 1$ on timeout
 - $cwnd = cwnd/2$ on triple dup ACK
- ❖ TCP-newReno
 - TCP-Reno + improved fast recovery (SKIPPED)
- ❖ TCP-SACK (NOT COVERED IN THE COURSE)
 - incorporates selective acknowledgements



Quiz: TCP Congestion Control?



In the figure how many congestion avoidance intervals can you identify?

- A. 0
- B. 1
- C. 2
- D. 3
- E. 4

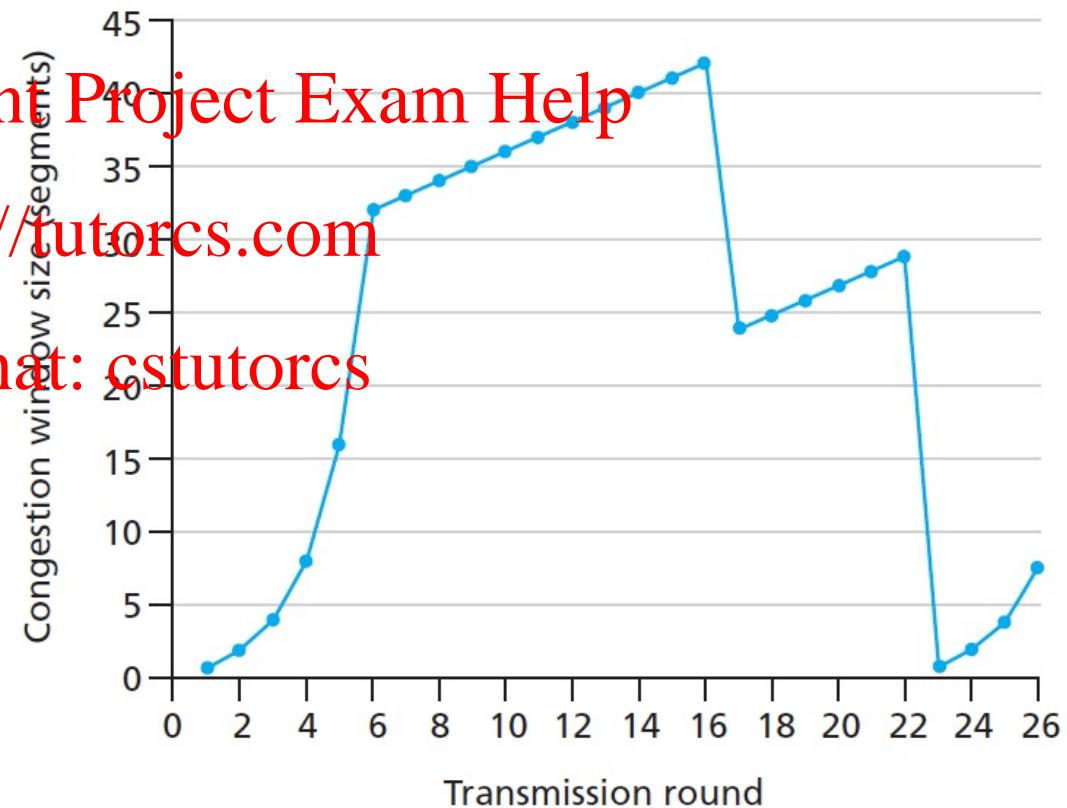
ANSWER: C

Assignment

Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs



Quiz: TCP Congestion Control?



In the figure how many slow start intervals can you identify?

- A. 0
- B. 1
- C. 2
- D. 3
- E. 4

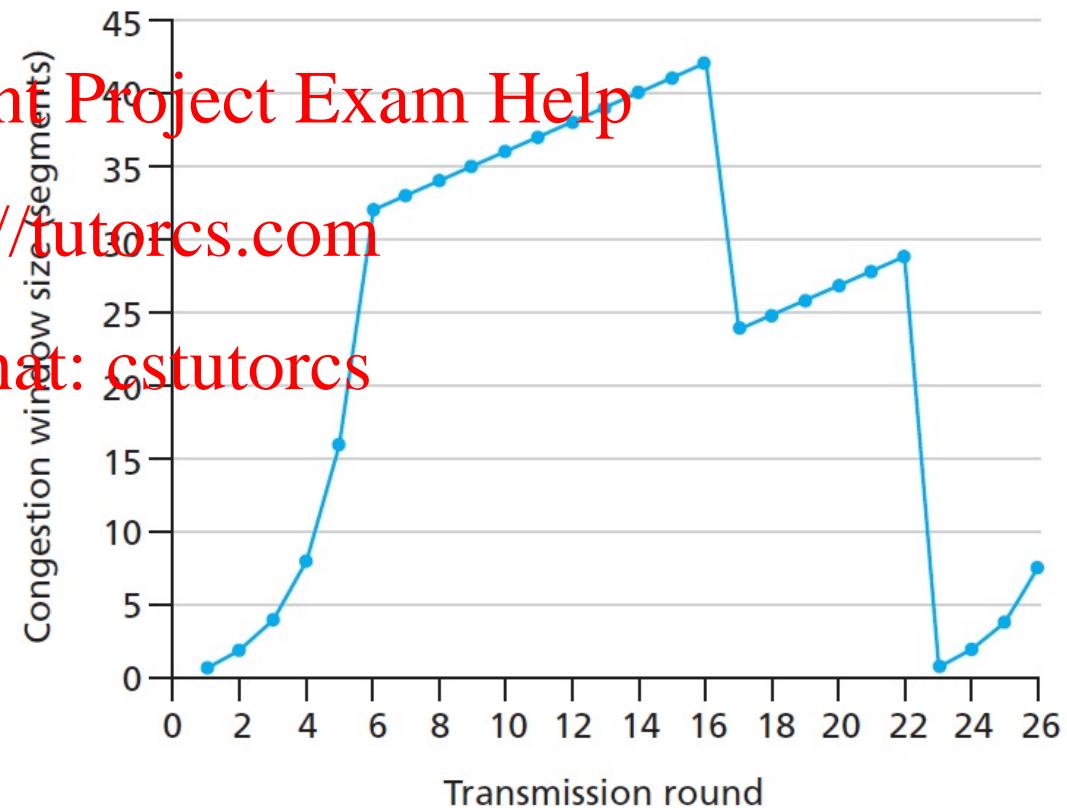
ANSWER: C

Assignment

Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs



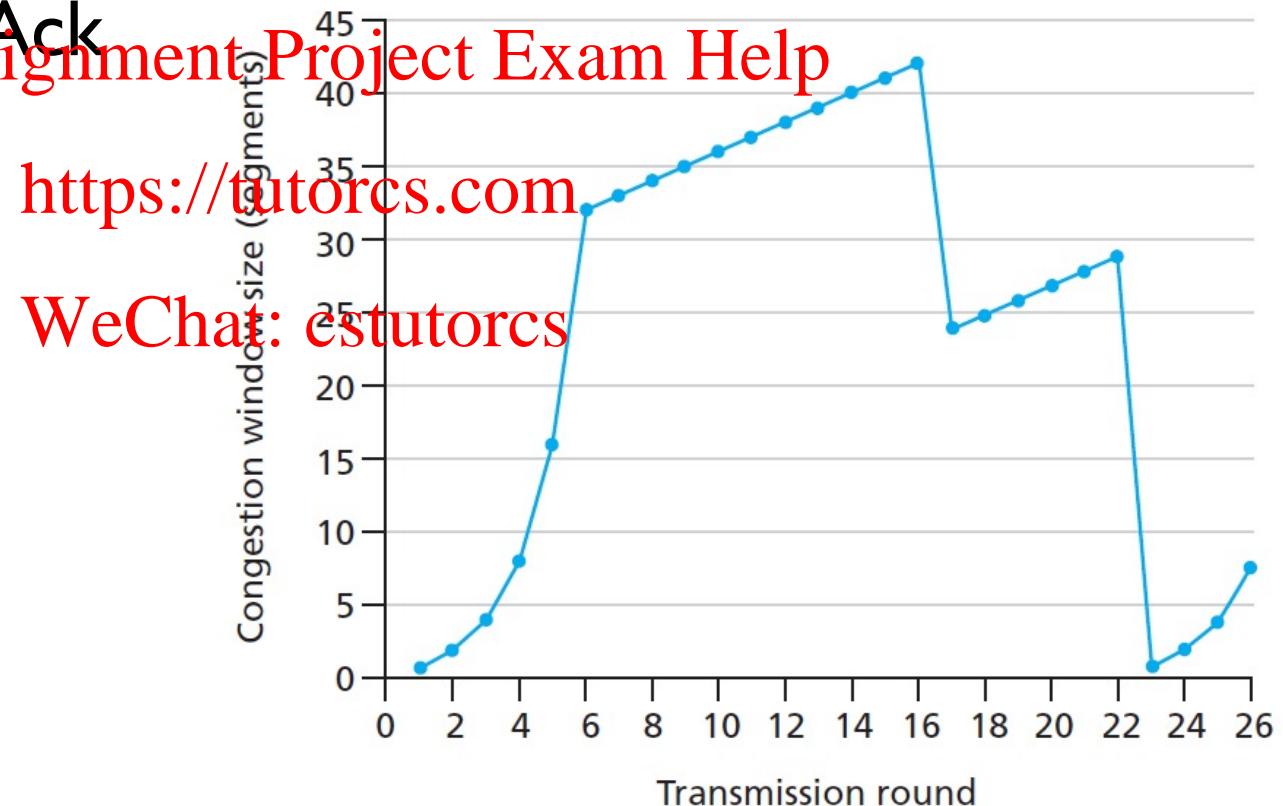
Quiz: TCP Congestion Control?



In the figure after the 16th transmission round, segment loss is detected by _____ ?

- A. Triple Dup Ack
- B. Timeout

ANSWER: A



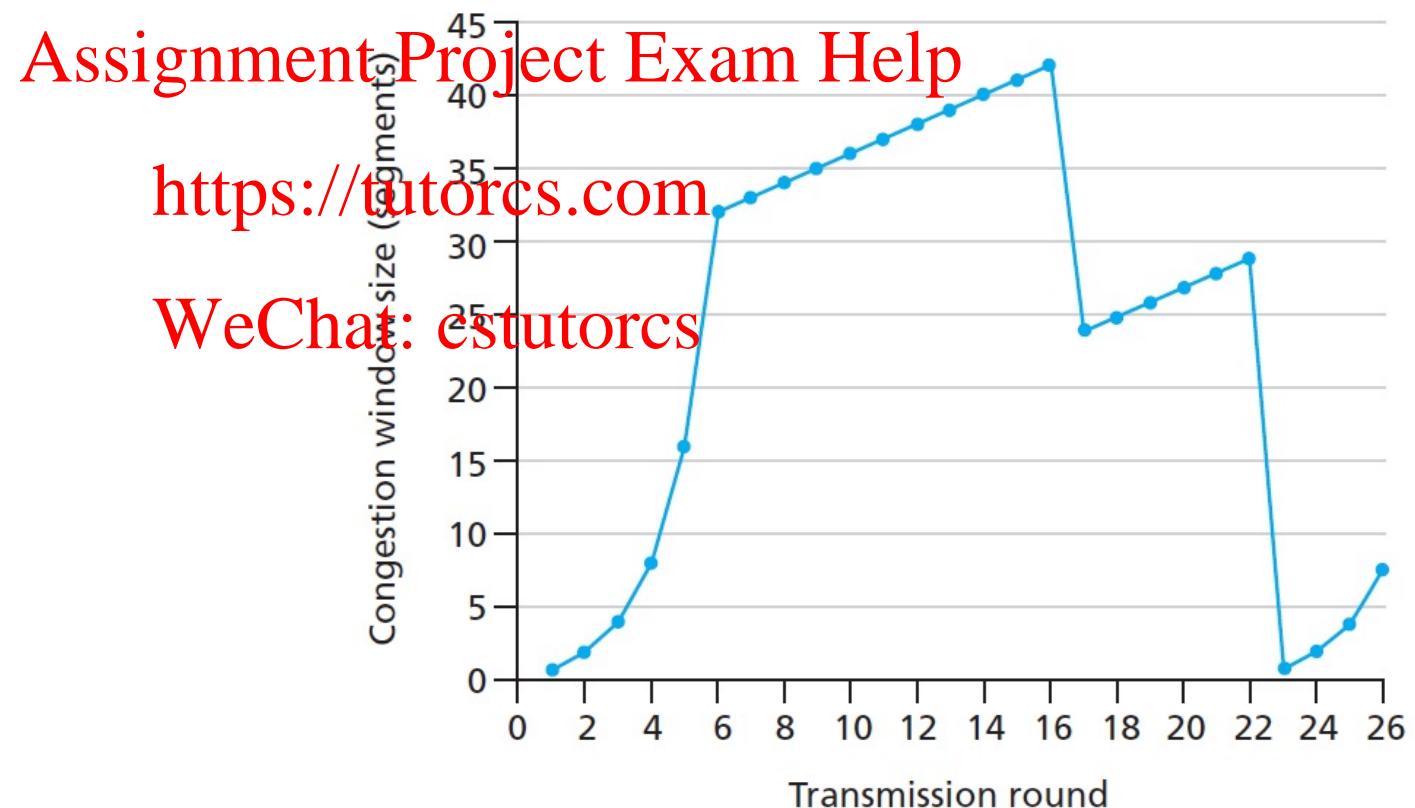
Quiz: TCP Congestion Control?



In the figure what is the initial value of sstresh (steady state threshold)?

- A. 0
- B. 28
- C. 32
- D. 42
- E. 64

ANSWER: C



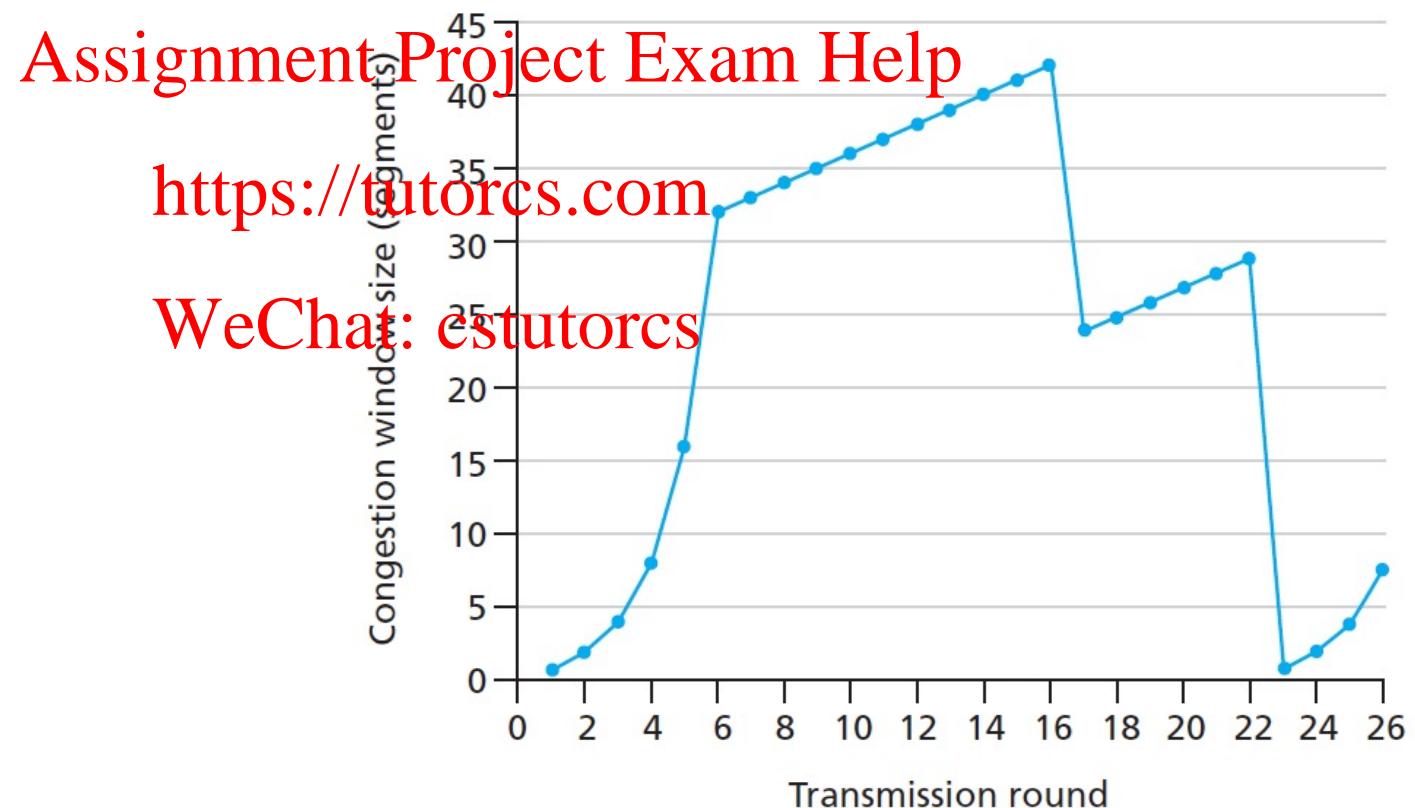
Quiz: TCP Congestion Control?



In the figure what is the value of ssthresh (steady state threshold) at the 18th round?

- A. 1
- B. 32
- C. 42
- D. 21
- E. 20

ANSWER: D



Transport Layer: Summary

- ❖ principles behind transport layer services:
 - multiplexing, demultiplexing
 - reliable data transfer
 - flow control
 - congestion control
 - ❖ instantiation, implementation in the Internet
 - UDP
 - TCP
- Assignment Project Exam Help
<https://tutorcs.com>
WeChat: cstutorcs
- next

- ❖ leaving the network “edge” (application, transport layers)
- ❖ into the network “core”