

SENDING & RECEIVING BUFFERS

TCP $\frac{2}{16}$

One implementation \rightarrow use circular array of 1-byte locations
Normally buffers are hundreds or thousands of bytes.

At sender's side 3 types of chambers

- a: Empty chambers that can be filled by sender process
- b: Bytes have been sent but not yet acknowledged.
- c: TCP keeps these bytes in buffer until ack is received.
- c: Bytes to be sent by the sender TCP.

At receiver side, circular buffer has 2 areas

- a: Empty chambers to be filled by bytes from the N/W.
- b: Chambers containing received bytes, that can be read by receiving process.

When a byte is read by receiving process,
chamber is recycled & added to pool of empty chambers.

SEGMENTS

TCP groups a no of bytes together into a packet called a segment. TCP adds header to each segment. Segments are encapsulated in IP datagram & transmitted.

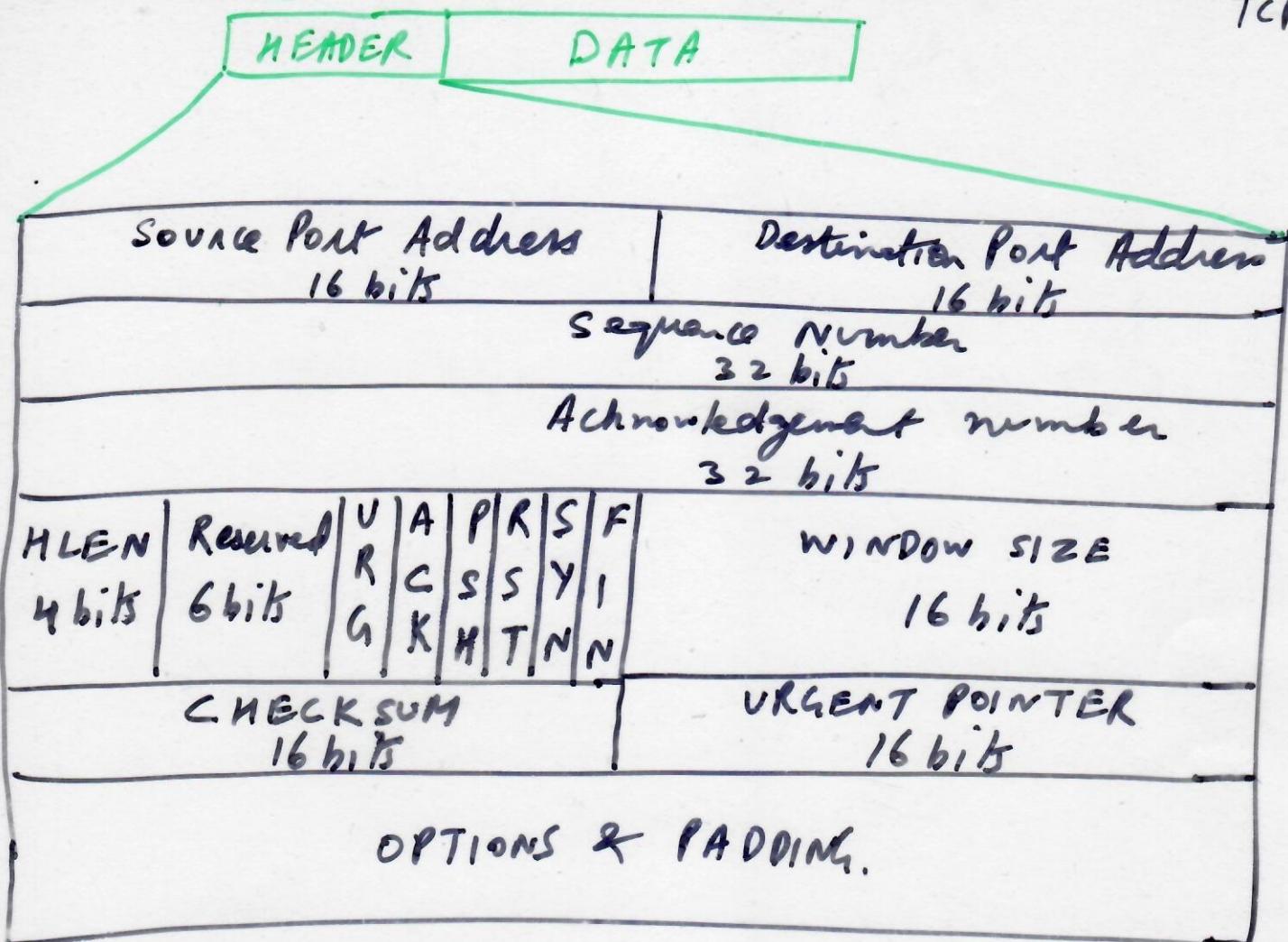
Segments are not necessarily of same size & are of the order of hundreds of bytes, if not thousands of bytes.

TCP offers FULL DUPLEX service.

TCP FEATURES TCP uses sequence no & acknowledgement no to keep track of bytes number in the segments

TCP nos the bytes stored in buffers with random nos generated between $0 - 2^{32} - 1$.

Sequence no is given to each segment that is being sent. Seq. No = No of 1st byte in that segment.

SEGMENT

Header = 20 - 60 bytes

Header = 20 bytes, then no options.

Source Port NO, Dest Port NO, Sequence NO = 1st byte in the segment i.e. data segment.

Ack NO, Header Length is between 20 - 60 bytes.

Value of this field between 5 ($5 \times 4 = 20$) & 15 ($15 \times 4 = 60$).

Reserved — for future use.

Control defines 6 different control bits or flags.

One or more bits can be set at a time.

URG : Urgent pointer is valid : used for urgent data.

ACK : Ack is valid

PUSH : Push the data

RESET : Connection must be reset.

SYN Synchronizes seq. nos during connection. TCP 4
FIN Terminates the connection. TCP 16

WINDOW SIZE Max size 65535 bytes. Window size is generally called receiving window (rwnd) & is determined by the receiver.

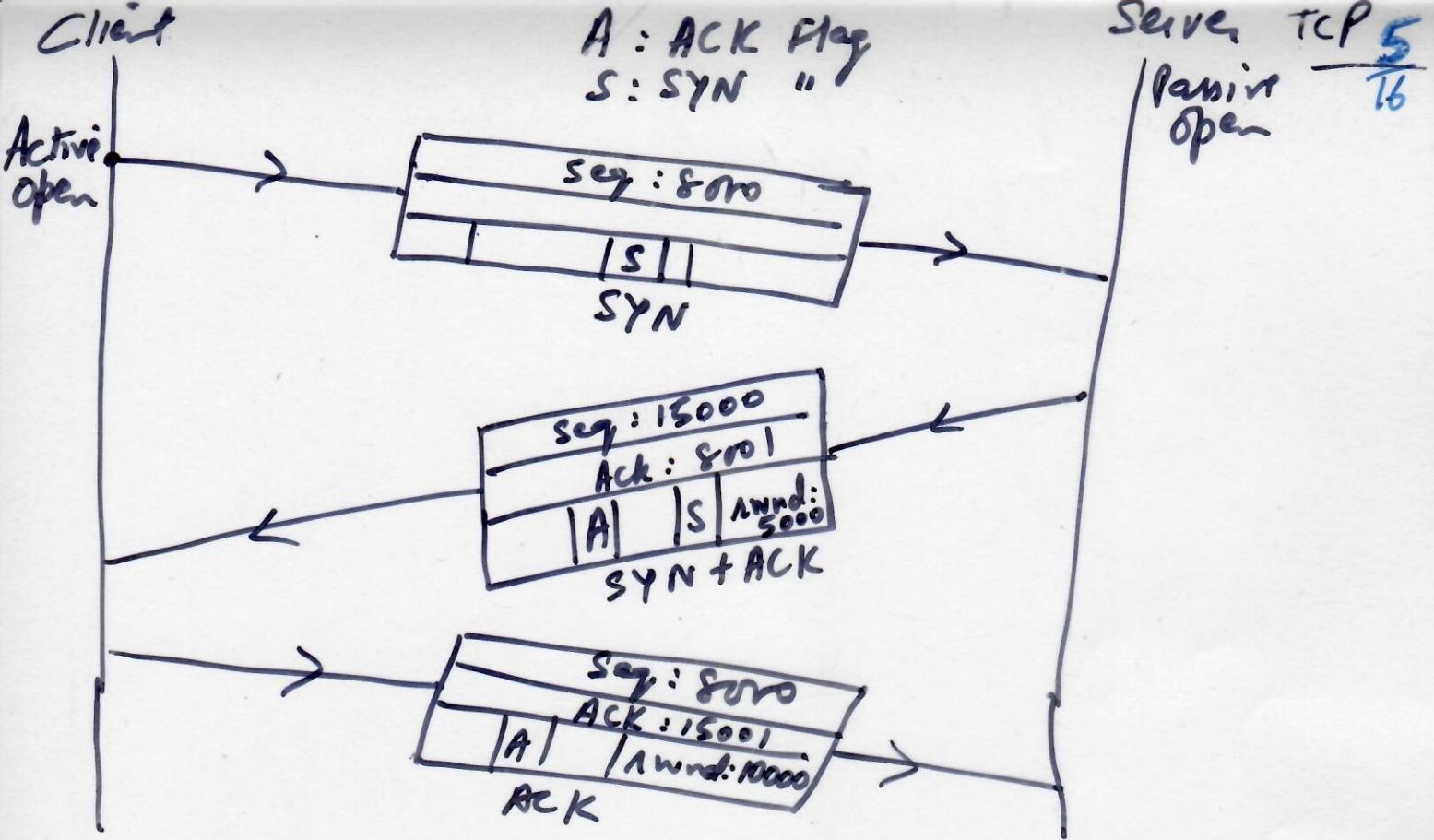
CHECKSUM

OPTIONS can be upto 40 bytes of information.

TCP CONNECTION establishes a virtual connection using single virtual pathway for entire message facilitates ack process as also retransmission of damaged or lost frames. TCP controls the connection.

Connection Establishment (Three Way Handshaking)

- a. Process starts with server.
- b. Server program tells its TCP that it is ready to accept a connection. This is request for passive open.
- c. Client program issues a request for an active open.
- i. Client sends 1st segment, SYN segment, in which only SYN flag is set. Client sends initial sequence no (ISN) (ISN). This segment doesn't have ACK no. SYN segment = control segment, so no data.
SYN " " has a sequence no
- ii. Data transfer sequence no is one more than SYN segment seq. no.



2. Server sends SYN+ACK segment. Serves 2 purpose
 It is SYN segment for communication in other direction.
 Server uses " " to initially assign a seq no for numbering
 the bytes sent from server to client.
 Server ACKnowledges.

3. Client sends ACK segment, just an ACK segment.
 Doesn't consume any seq no.
 Some implementations carry user data, where seq no is consumed.

Issues Involved in Connection Establishment

(a) Simultaneous Open When both processes issue active open..

Both TCPs send SYN+ACK segments to each other & one single connection is established between them.

(b) SYN Flooding Attack When a malicious attacker sends large SYN segments to server pretending each config from different IP addresses.

76

Detecting corrupted segments, lost segments, out of order segments & duplicated segments.

Error correction & detection by checksum, ACK, Time-out.

CHECKSUM

Corrupted segment dropped by dest TCP is considered lost. TCP uses 16 bit checksum.

ACKNOWLEDGEMENT

Some Imp rules

1. Piggybacking, reduces no of segments.
2. Receiver needs to delay sending an ACK segment if there is only one outstanding in-order segment.
3. When a segment arrives with a seq no that is expected by the receiver & the previous in-order segment has not been acknowledged, the receiver immediately sends an ACK segment.
4. When a segment arrives with an out-of-order seq no that is higher than expected, receiver immediately sends ACK segment announcing the seq no of the next expected segment.
5. When a missing segment arrives, the receiver sends an ACK segment to announce the next seq no expected. This informs receiver that segments reported missing have been recd.
6. If a duplicate segment arrives, receiver sends ACK immediately. This solves some problems when an ACK itself is lost.

ACK Type

- 1: Accumulative ACK Receiver advertises next byte it expects to receive, ignoring all segments recd out of order. Discarded, but on duplicate segments are not reported. 32 bit ACK field in TCP header is used for accumulative ACK & its value is valid only when ACK flag bit = 1
- 2: Selective ACK doesn't replace ACK, but reports additional info to sender. Out of order & duplicated segments are reported. SACK is implemented as option at end of TCP header.

RETRANSMISSION (a) Occurs if retransmission timer expires.

(b) Sender receives 3 duplicate ACKs.

Retransmission After RTO (Retransmission Time - Out)

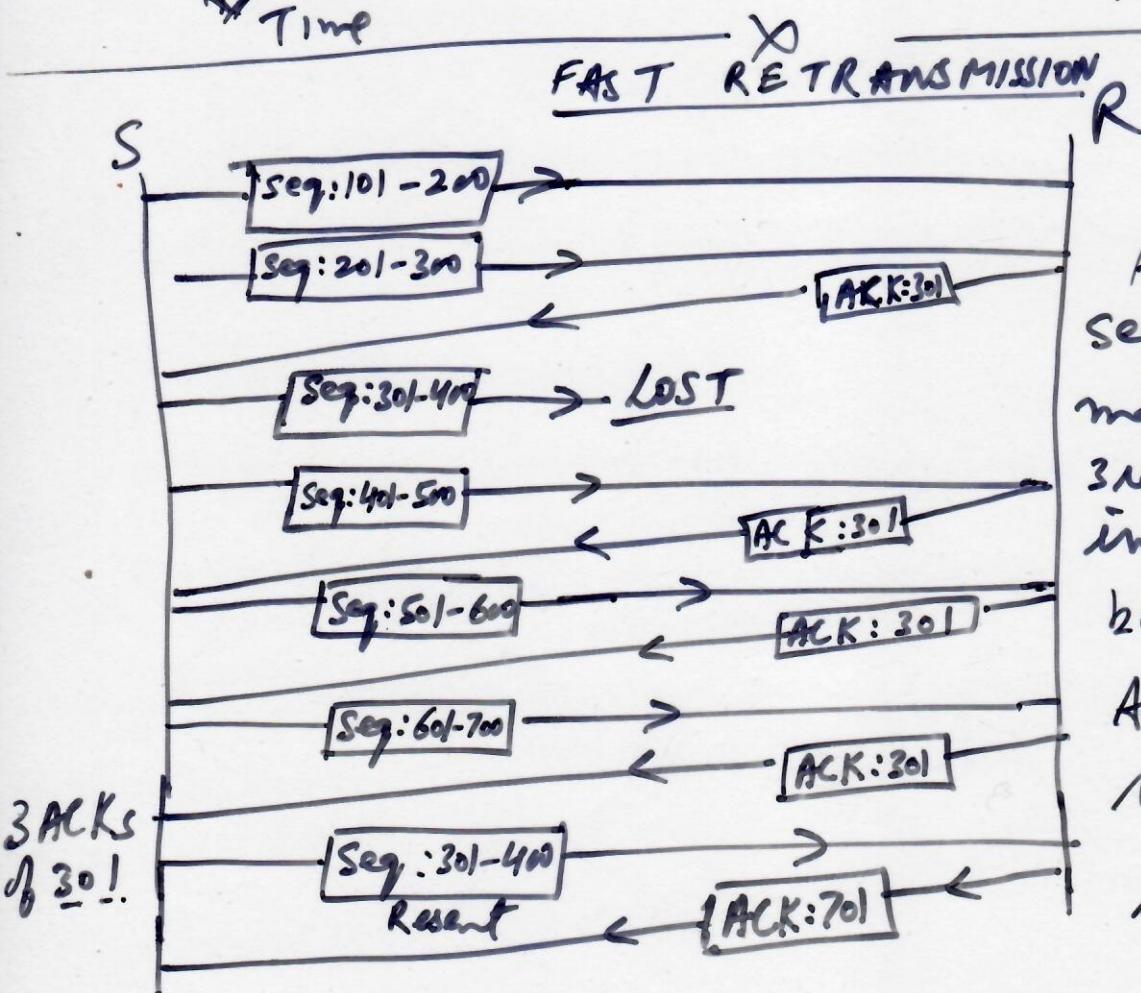
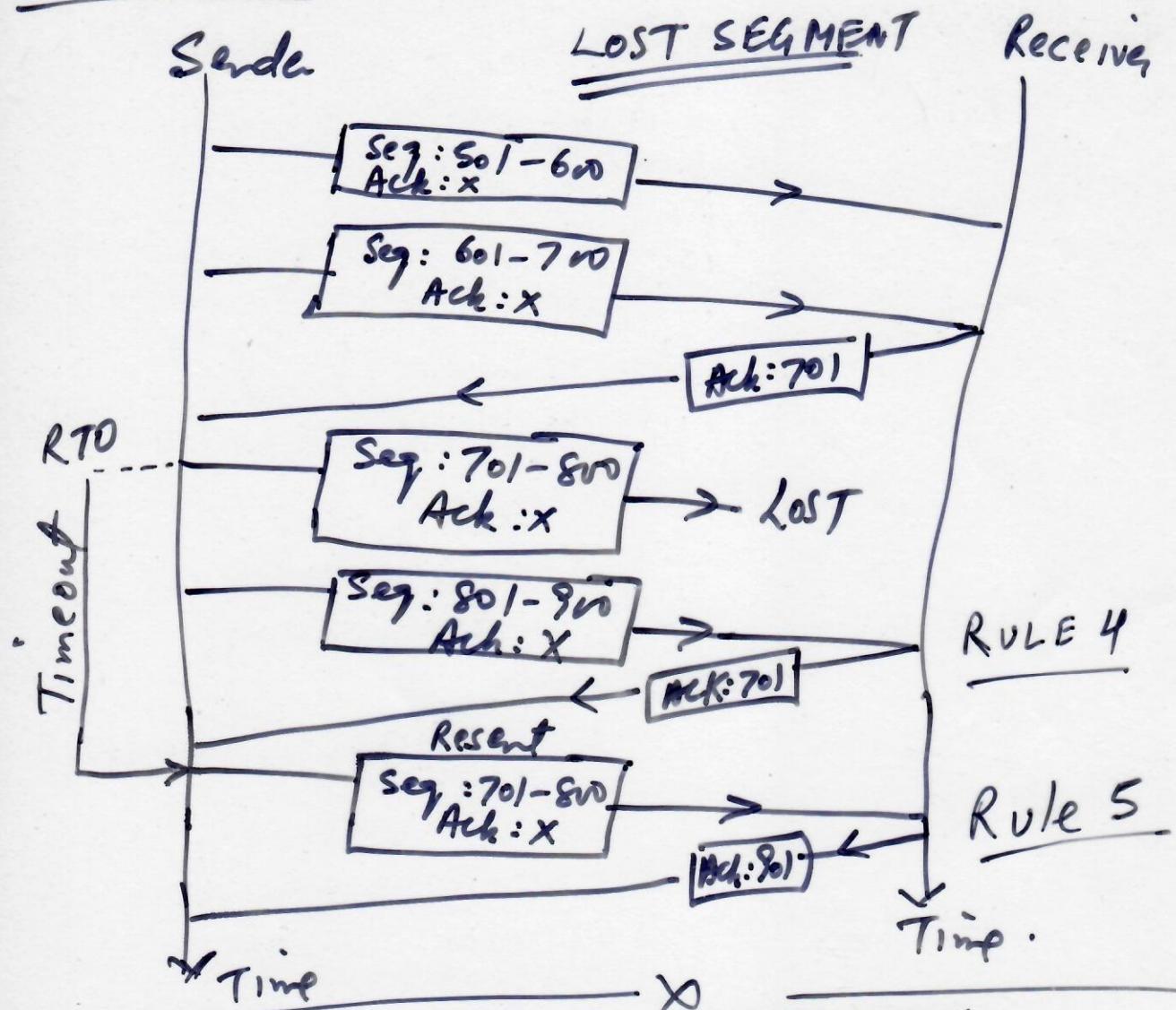
Source TCP starts an retransmission time out (RTO) timer for each segment sent.

RTO is dynamic & is updated based on round trip time (RTT) of segments.

Out of Order Segments Originally TCP discarded all out of order segments. Now, out of order segments are not discarded but stored temporarily till missing segments arrive, then given to process (in order).

Some Scenarios

TCP 8
16

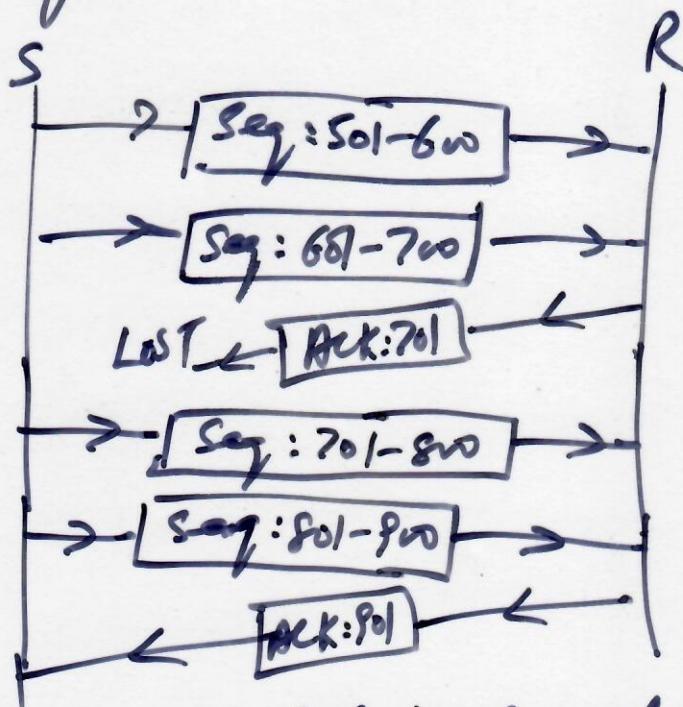


Delayed Segment Diff. write leads to delayed segments. " may arrive after it has been resent. (Duplicate segment). TCP 8-9
16

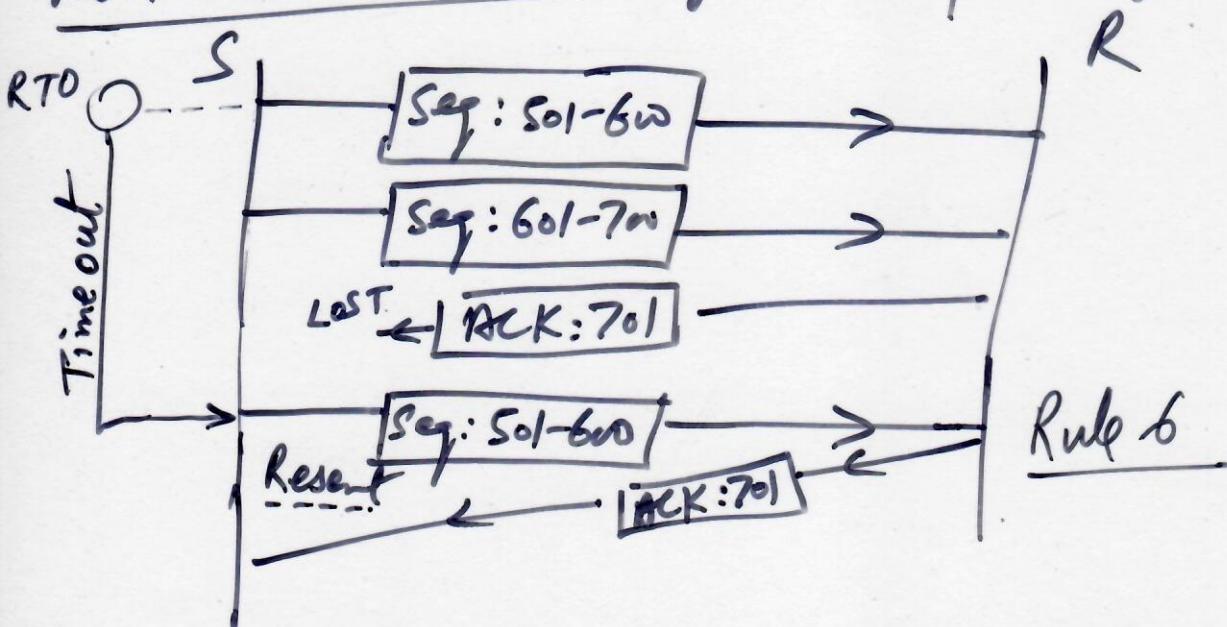
Duplicate Segment Duplicate seg is handled as follows:

when a seg arrives that contains a seq no less than the previous by Acknowledged bytes, it is discarded.

Automatically corrected Lost ACK



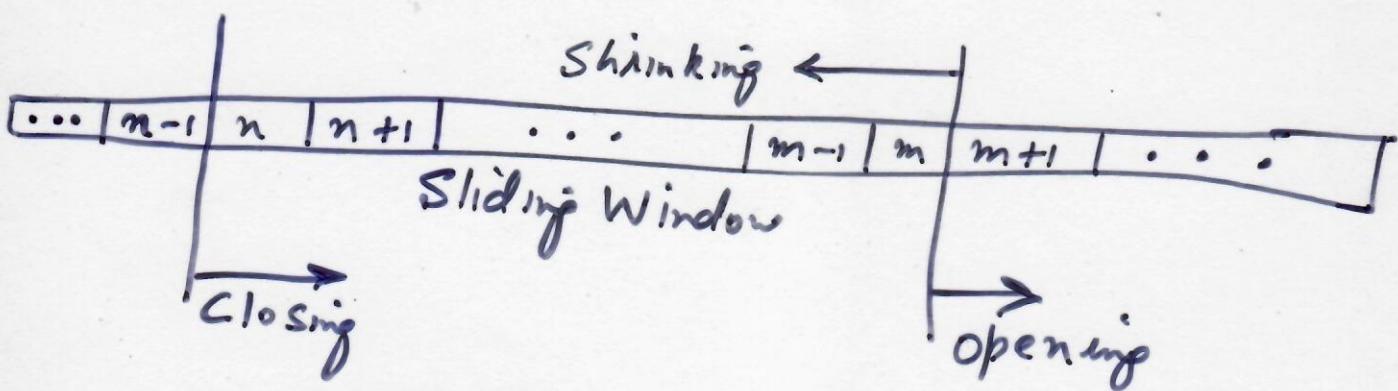
Lost ACK corrected by Resend ip a seg



FLOW CONTROL regulates the amount of data a source can send before receiving an acknowledgement from the destination. TCP defines a window that is imposed on the buffer of data delivered from the application program & ready to be sent. TCP sends an amount of data defined by the sliding window protocol.

Sliding Windows Protocol TCP uses sliding windows protocol for flow control. In this method, host uses a window for outbound communication. (sending data). Window is portion of the buffer containing bytes received from the process.

Bytes inside the window are the bytes that can be in transit. The imaginary window has 2 walls, one left & one right. The window is called sliding window because the right & left walls can slide as shown below:-



$$\text{Window Size} = \min(\text{Awd}, \text{Cwd})$$

Window is opened, closed or shrunk. These are in the control of the receiver, not the sender.

Opening window implies moving the right wall to the right. This allows new bytes in the buffer that are eligible for sending. Closing windows is moving left wall to the right which means some bytes have been acknowledged.

Shrinking is moving ^{right} left wall to the ~~right~~ left. This is strongly discouraged & not allowed in some implementations, because it means revoking the eligibility of some bytes for sending. If sender has already sent these bytes, then previous ACKs have to be rechecked, which is a problematic area.

However there is an exception i.e. the receiver can temporarily shutdown the window by sending $1 \text{ wnd} = 0$. This implies receiver doesn't want data from sender for a while.

Some facts about TCP's sliding windows

1. Size of the window is the lesser of 1 wnd & $c \text{ wnd}$.
2. Source doesn't have to send a full window's worth of data.
3. Window can be opened or closed by the receiver but should not be shrunk.
4. Destination can send an ACK at any time as long as it doesn't result in a shrinking window.
5. The receiver can temporarily shut down the window. The sender, however, can always send a segment of one byte after the window is shut down.

Silly Windows Syndrome Problems can occur in sliding window operation when sender creates data slowly or receiver consumes data slowly or both. Any of these can result in reduced efficiency. For e.g. if TCP sends segments of 1 byte of data, it implies 41 byte datagram (20 bytes of TCP header + 20 bytes of IP header). This problem is called silly windows syndrome.

Syndrome created by Sender The process writes 1 byte at a time into the buffer of the sending TCP.

Solution is to prevent the sending TCP from sending data byte by byte - TCP must be forced to wait for more data. How long should TCP wait? If longer is wait it delays the process. If wait is small, it may result in small segments. Nagle found elegant solution.

Nagle's Algorithm

1. The sending TCP sends the 1st piece of data it receives from the sending application program even if it is 1 byte.
2. After sending 1st segment, data is gathered in the OIP buffer & waits until ACK is received or enough data is received to fill a maximum-size segment. Now TCP can send the segment.
3. Step 2 is repeated for the rest of transmission.

Nagle's Algo takes into account speed of application program & the speed of N/w that transports the data.

If application program is faster than the N/w, segments are larger.
" " " " " slower " " " , " " " smaller.

Syndrome created by the Receiver

Receiving TCP may create a silly window syndrome if it is serving an application program that consumes data slowly, for e.g., 1 byte at a time. Suppose sender creates data in blocks of 1 kbyte, but receiver consumes data 1 byte at a time. Suppose TCP buffer of receiver is 4 kbytes. Sender sends the 1st 4 kbytes of data. Receiver stores it in its buffer, which is full. Receiver advertises window size of zero, which implies sender should stop sending. Receiver reads the 1st byte of data from TCP buffer. Receiver advertises window size of 1 byte & sender sends 1 byte. The procedure will continue.

Two solutions have been proposed

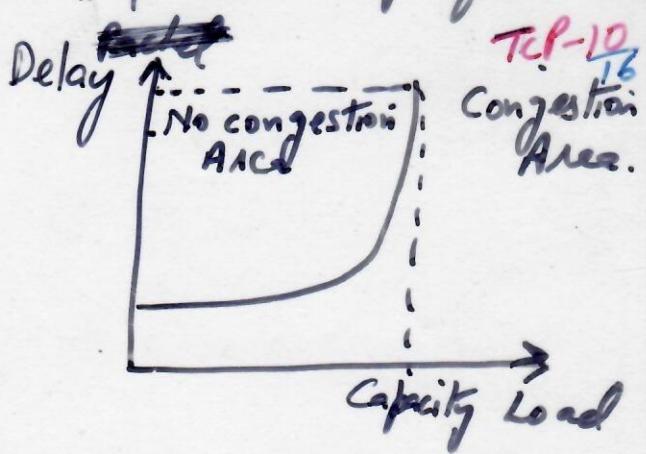
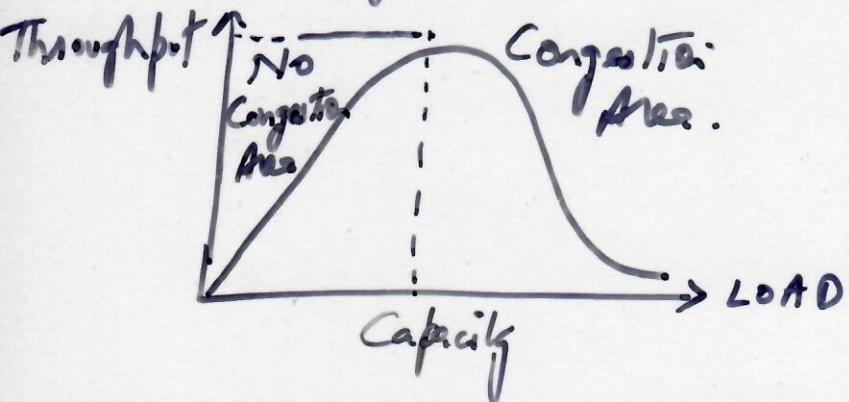
Clark's Solution: To send an ACK as soon as data arrives, but to announce a window size of zero until there is enough space to accommodate a segment of maximum size or until half of the buffer is empty.

Delayed ACK: When segment arrives, it is not ACK immediately. Receiver waits until there is a decent amount of space in its incoming buffer before ACKing the arrived segments. Delayed ACK prevents the sender TCP from sliding its window. After it has sent the data in the window, it stops. This kills the syndrome.

Delayed ACK has the advantage of reduced traffic. Receiver doesn't have to acknowledge each segment. Disadvantage of delayed ACK is, it may force the sender to retransmit the unacknowledged segments.

The protocol balances the advantages & disadvantages. It now defines that the ACK should not be delayed by more than 500 ms.

CONGESTION CONTROL involves 2 factors that measure performance of a NW: delay & throughput.



CONGESTION CONTROL MECHANISMS

(a) Open Loop Congestion Control policies are applied to prevent congestion. Congestion control by — source or dest.

Brief List of Policies for congestion prevention are:-

- Retransmission Policy Good retransmission policy can prevent congestion.
- Acknowledgment Policy imposed by receiver may also affect congestion.
- Discard Policy Good discard policy possessed by routers may prevent congestion & at same time will not harm the integrity of transmission.

(b) Closed Loop Congestion Control policies are applied to remove congestion after it happens. Few mechanisms are

- Back pressure when route is congested, it can inform the previous router to reduce rate of outgoing packets. This action can be recursive upto router just prior to source.
- Choke Point packet sent by router to source to inform congestion.
- Implicit Signaling source can detect implicit signal without warning of congestion & slow down its sending rate.
- Explicit Signaling Routers that experience congestion can send an explicit signal, setting of a bit in packet to inform sender/receiver of congestion.

CONGESTION CONTROL IN TCP

$CWnd = \text{congestion window size}$ TCP-11

Congestion Window Actual window size = minimum (1 wnd , $CWnd$) 16

Congestion Policy TCP general policy for handling congestion is based on 3 phases :- slow start, congestion avoidance, congestion detection.

In slow start, sender starts with very slow rate, but increases rate rapidly to reach threshold. After reaching threshold data rate is reduced to avoid congestion. Finally if congestion is detected, sender goes back to slow start or congestion avoidance phase based upon how congestion is detected.

Slow Start starts with $CWnd = 1$, then increased exponentially upto slow start threshold. Then slow start stops & congestion avoidance phase starts. $CWnd = 1$ implies 1 MSS

Slowstart threshold = 65535 bytes (usually) (Max segment size)

Congestion Avoidance When threshold is reached, the window size is increased additively, until congestion is detected.

Congestion Detection (a) Retransmission Timer expires, then .

Threshold window size = $\frac{1}{2} \times \text{Current window size}$

& Slow start phase starts where $CWnd = 1$

(b) When 3 ACKs are received, then

Threshold window size = $\frac{1}{2} \times \text{Current window size}$

$CWnd = \text{slow start threshold}$

& Congestion Avoidance phase starts.

TCP Timers

- (a) Retransmission
- (b) Persistence
- (c) Keepalive
- (d) TIME-WAIT

Retransmission Timer To retransmit a lost segment, TCP TCP RTT RTT employs a " " that handles Retransmission Time-out (RTO), waiting time for an ACK of a segment.

for every segment Retransmission Timer is created.

2 situations may occur

- (a) If ACK is recd for a particular segment before timer expires, timer is destroyed.
- (b) If timer expires, before ACK arrives, segment is transmitted & timer reset.

ROUND TRIP TIME (RTT) To calculate RTO, we require RTT.

Measured RTT (RTTM) Only one RTTM at any time.

Smoothed RTT Fluctuation in RTTM is so high in today's internet, that it can't be used for retransmission time out purposes. Most implementations used smoothed RTT called RTTs which is weighted average of RTTM & previous RTTs.

Original \rightarrow No value

After 1st measurement \rightarrow $RTTs = RTTM$

After any other " \rightarrow $RTTs = (1-\alpha)RTTs + \alpha \cdot RTTM$

Value of α is implementation dependent, but normally $= \frac{1}{8}$.

RTT Deviation (RTTD) Most implementations use this.

Original \rightarrow No value

After 1st measurement \rightarrow $RTTD = RTTM / 2$

After any other " \rightarrow $RTTD = (1-\beta)RTTD + \beta / |RTTs - RTTM|$

β is implementation dependent, but usually = $\frac{1}{4}$.

TCP 13
16

Retransmission Time Out (RTO) based upon RTTs & its deviation.

Original \rightarrow Initial value

After any measurement $\rightarrow RTO = RTT_s + 4RTT_d$

KARN'S ALGORITHM If segment not acknowledged during retransmission period is retransmitted. Sender TCP receives ACK for this segment, it doesn't know if ACK is for original segment or for the retransmitted one. The value of new RTT is based upon departure of segment. However, if original segment was lost & ACK is for the retransmitted one, the value of current RTT must be calculated from the time the segment was retransmitted. This is a dilemma which was solved by Karn. Karn solution is simple. Don't consider RTT of a retransmitted in calculation of new RTT. Do not update value of RTT until you send a segment & receive an ACK w/o the need for retransmission.

Exponential Back Off What is value of RTO if retransmission occurs?

Value of RTO is doubled for each retransmission. If segment is retransmitted once, value is 2 times the RTO, if it transmitted twice, value is four times the RTO & so on.

Persistence Timer Receiving TCP announces window size = 0. Sending TCP stops transmitting until receives an answer non-zero window size. This ACK can be lost leading to deadlock. There is no retransmission timer for a segment containing only an ACK.

To resolve deadlock TCP uses Persistence Time. Sender TCP starts persistence timer when it receives an ACK with window size = 0.

When time elapses, sending TCP sends special segment called a probe, containing just one byte of data. It has a seq no, but this seq no is never acknowledged. Probe alerts receiver that ACK was lost & must be resent.

Value of Persistence Time = Value of Retransmission Time.

If no respond from receiver, another probe segment & value of persistence timer is doubled & reset.

This process of sending probe segment is repeated till value of persistence timer reaches threshold (usually 60s).

After that sender sends one probe segment every 60 s until window is opened.

Keepalive Time: To prevent long idle connection between 2 TCPs. Suppose client opens a TCP connection & becomes silent, connection opens forever.

Most implementations equip server with a keepalive timer. Each time server hears from client, it resets timer. Time out is usually 2 hours. If server doesn't hear from client after 2 hrs, it sends probe segment.

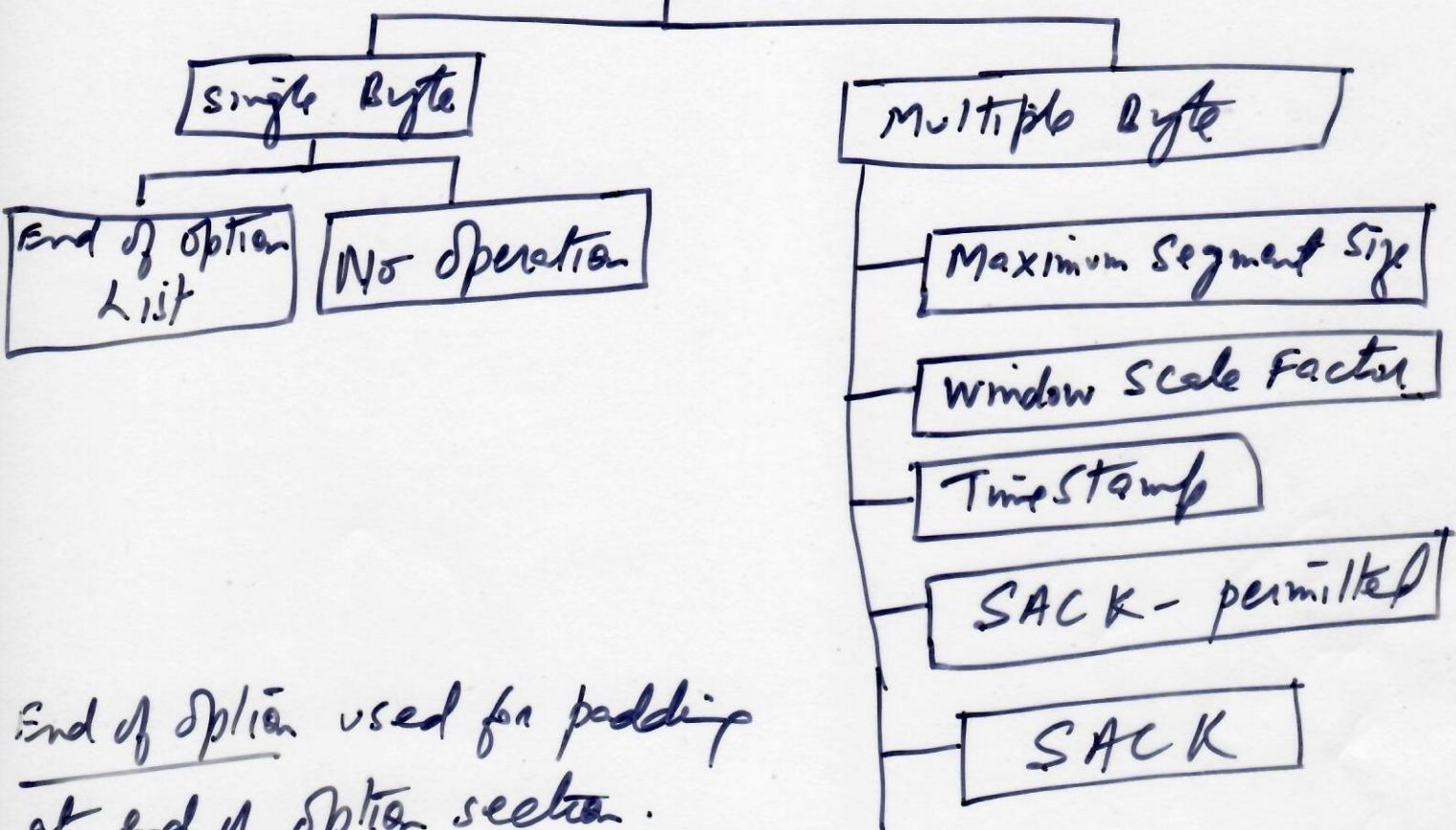
If no response after 10 probes, each of which is 755 apart, it assumes client is down & terminates connection.

Time Wait Timer used during connection terminated.

OPTIONS convey additional info to destination or align other options.

TCP IS
16

Options



End of Option used for padding at end of option section.

After this option receiver looks for payload data.

No operation used as a filler. However, it normally comes before another option to make an option fit in a 4 word slot.

Maximum Segment Size (MSS) defines largest unit of data that can be received by dest TCP. It defines max data size & not max size of segment. $MSS = 16 \text{ bit field}$. Default value of MSS = 536 bytes.

Window Scale Factor 16 bit field, defines size of sliding window, which can range from 0 - 65535 bytes. This size may not be sufficient. To increase window size, window scale factor is used.

new window size = window size in header $\times 2^{\text{window scale factor}}$
Although scale factor could be as large as 255, the largest value allowed by TCP/IP is 14 i.e. max window size is $2^{16} \times 2^{14} = 2^{30}$ which is less than max value of seq no.

Timestamp 10 byte option. Timestamp has 2 applications
(a) measures RTT (b) Prevents wrap around seq. nos.

Last
Slide
TCP

Identification of segment = combination of seq no + timestamp

SACK Permitted and SACK Options

SACK Permitted \Rightarrow 2 bytes. SACK permitted used only during connection establishment. Host that sends SYN segment adds this option to show that it can support SACK option.

SACK option, of variable length used during data xfer only if both ends agree. The option includes list of blocks arriving out-of-order.

