

NAME: Bhukya Vasantha Kumar
ROLL: B180411CS (A-Batch)

= "Module 3 Notes" =

Transmission Control Protocol (TCP)

It is one of the main protocols of the Internet Protocol (IP) suite. The suite is also known as TCP/IP since it originated with initial implementation of IP.

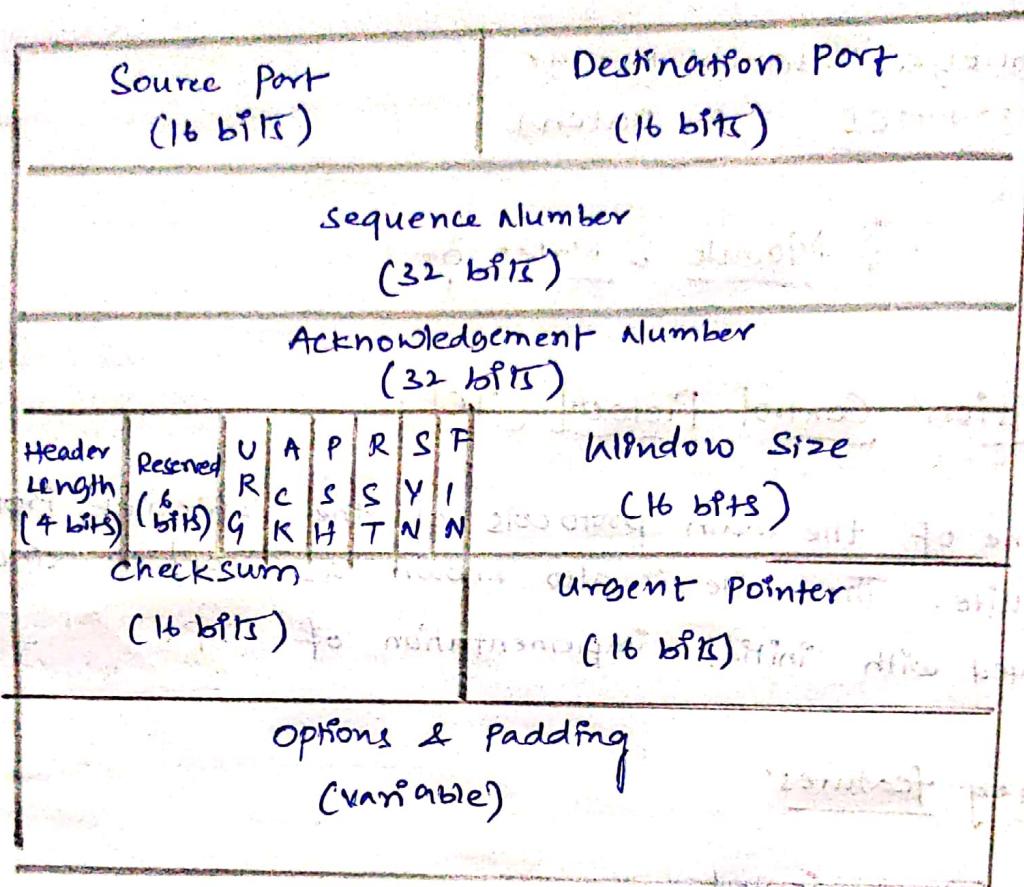
Some key features:

- Connection oriented
- Point-to-point communication between 2 end points.
- Reliable data transfer via packetization with sequence numbers, header, checksum, etc.
- Full duplex communication : bidirectional & symmetric.
- Reliable connection startup
- Graceful connection shutdown
- Reliable & adaptive retransmission

Other perks:

- End-to-End service: Application handle their own features at the end points and intermediary nodes only relay the payload.
- Reliable Transmission: If any packet loss occurs, the packets are retransmitted. The 3 way handshake ensures a connection before a payload is sent.

Header Format



Source port: identifies source user process

Destination port: identifies destination user process

Sequence number: byte number of the first byte in segment.

Acknowledgment number: byte number receiver expects next.

Header length: indicates the length of TCP header.

Control Flags: 1 bit control bit that control connection establishment, termination, flow control, etc.

URG: Urgent pointer is valid

ACK: Acknowledgment number is valid

PSH: Request for push

RST: Reset Connection

SYN: Synchronize sequence number

FIN: Terminate connection

Window size: Indicates window size for flow control

Checksum: holds checksum for error control

Urgent pointer : If URG is set, then this 16-bit field is an offset from sequence number indicating last urgent byte. (2)

options: optional field can be 0-320 bit

TCP Checksum

- It is a 16 bit one's complement of the one's complement sum of a pseudo header of information from the IP header, the TCP header and the data padded with zero octets at end to make a multiple of 2 octets.
- Filled with zeroes initially
- TCP length is not transmitted but used in calculations.

TCP Service Requests

- Unspecified passive open : Listen for connection requests from user.
- Full passive open : Listen for connection requests from specified user
- Active open : Request connection
- Active open with data : Request connection & transmit data
- Send : send data
- Allocate : Issue incremental allocation for receive data
- Close : close connection gracefully
- Abort : close connection abruptly
- Status : Report connection status

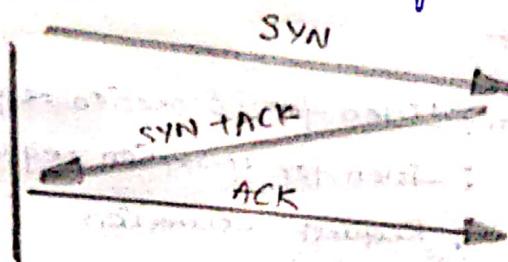
TCP Service Responses

- open ID : Informs name assigned to pending request
- open failure : Open request failed
- open success : Open request succeeded
- Deliver : Reports arrival of data
- closing : Remote TCP issued a close request
- Terminate : Connection terminated
- status response : Reports connection status
- service response : Reports service requests / internal error
- Error

TCP Mechanisms

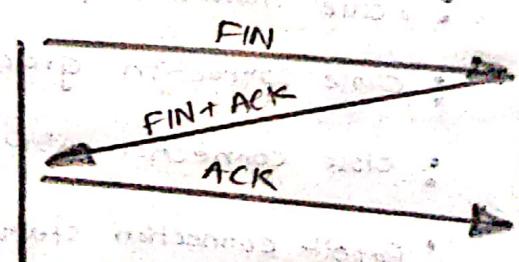
- connection Establishment

3 way handshake with SYN flag set → request for connection



- connection Termination

close with FIN flag set



Data Transfer

- stream: Every byte numbered modulo 2^{32}
- Header contains sequence number of 1st byte
- Flow control: credit = no of bytes
- Data transmitted in intervals determined by TCP.
- Push: send now

- urgent: send data in ordinary data stream with urg.
- If TDPV not intended for this connection received → "reset" flag set in outgoing segment.

Implementation policies

- Send policy

Too little: more overhead

Too large: delay

Push: send now

- Delivery policy

May store/deliver each in-order segment

Push: send now

- Accept policy

May/may not discard out-of-order segments.

- Retransmit policy

First only

Retransmit all

Retransmit individual

(Maintain separate timer for each segment)

- Ack policy

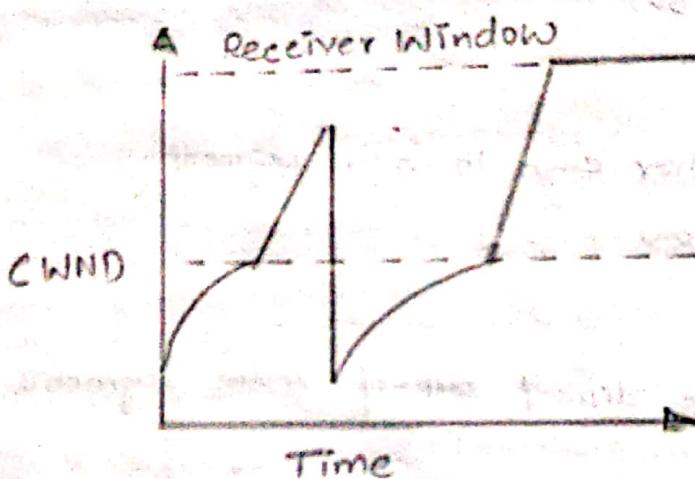
Immediate (No piggybacking)

Cumulative (Wait for outgoing data) time out

Slow start flow control

- Window → flow control avoids receiver overrun
- Needs congestion control → avoid network overrun
- Sender maintains 2 windows:
 - credits from receiver
 - congestion window from network
- Congestion window < receiver window

- Start with congestion windows of 1 segment → do not disturb existing connections too much.
- Increase Cwnd by 1 for every ACK received.
 - If packets lost → slow start threshold $\rightarrow Cwnd/2$
 - Set Cwnd = 1
 - Increment 1 per ACK until SS threshold
 - Increment $1/Cwnd$ per ACK afterwards



Fast Retransmit & Recovery

- If same packet acknowledgment (3 Nacks) \rightarrow
 - assume next packet lost
 - retransmit of right away.
 - retransmit only 1 packet
- useful when 1 packet lost, not multiple.

Selective ACK. (SACK)

- Initial Negotiation: sender to receiver "SACK permitted"
- Selective ACK: variable length - receiver to sender
- Left edge: 1st seq. no. in this block
- Right edge: Req. no. immediately after last seq. number in block.

- ACK field same : next byte receiver expects
- Missing segment received : ACK field advanced
- Receiver can send SACK only if sender has "SACK permitted" option in SYN segment of connection
- Data receiver can discard sacked data, sender must not until acknowledged.

Window Scaling Option

- Long fat pipe Networks (LFN) : satellite links

- Need very large window sizes

- Normally,
max window = $2^{16} = 64\text{K bytes}$

- Window scale option,

$$W = W * 2^{\text{scale}}$$

i.e. W = W * 2^{\text{scale}}

- Max window = $2^{16} \times 2^{255}$

- option sent only in SYN & SYN+ACK segment

TCP/IP Tools

- netlookup
- ping
- finger
- traceroute
- whois
- netstat
- ipconfig /ifconfig
- dig / host
- tcpdump
- wget, etc.

TCP Extensions for high speed networks

H-TCP (TCP for high speed & Long distance networks)

Why H-TCP?

Future communication & computer networks

- high speed & long distance connectivity
- require transport layer enhancement
- Develop a systematic framework for modifying the basic TCP algorithm to render it suitable in a way of variety of network types.

The network new variant, H-TCP is:

- suitable for deployment in high speed & long distance networks as well as connectional networks.
- shown to be fair when deployed in homogeneous network.
- friendly when competing with conventional TCP sources.
- shown to rapidly respond to changes in available bandwidth & utilize link bandwidth efficiently.
- shown to behave as a conventional TCP variant when deployed on conventional network types.

How does it work?

α_i , additive increase parameter for the scene is small for conventional networks (for backward compatibility) & large in high speed & long distance networks (for short duration congestion epoch even with large pipe sizes)

- Modify the basic TCP paradigm by adjusting the rate α_i at which a source needs to insert packets into network to reflect prevailing network conditions.

→ Make Δ_i increase as function of time elapsed since last packet drop by i^{th} source. (09)

High speed mode source increase function $\rightarrow \alpha_{iH}(\Delta_i)$

Low speed mode source increase function $\rightarrow \alpha_{iL}$

$$\Delta_i = \begin{cases} \alpha_i^2 & \Delta_i \leq \Delta^2 \\ \alpha_{iH}(\Delta_i) & \Delta_i \geq \Delta^2 \end{cases}$$

Where,

Δ_i = Time elapsed since last congestion event by i^{th} src.

α_i^2 = Increase parameter for low speed regime.

$\alpha_{iH}(\Delta_i)$ = Increase function for high speed regime.

Δ^2 = Threshold for switching from low to high speed regimes.

MuiltCP (Schizophrenic TCP)

Originally proposed for differentiated services, to be used with a pricing scheme. TCP flow behaves like a collection of several virtual flows.

Algorithms Modified:

④ Slow Start

For each ACK, triple the 'ewnd' for first ' k ' RTTs, double after that. Congestion window at the end of ' k ' RTTs = $3k$.

• Sum of congestion window of N flows = $N2^k$.

• Choose k such that, $3k = N2^k$.

Linear Increase

Increase the cwnd by N cwnd for each incoming ACK.

Multiplicative Decrease

For each loss, reduce cwnd for only one virtual flow,
i.e. reduce cwnd by $(\frac{1}{N}) * 0.5 * \text{cwnd}$.

Downsides:

- Cannot modify / emulate timeouts
- N set by user or not controlled by end-to-end congestion situations and conditions.
- Relies on policing and pricing to prevent congestion collapse.
- More bursty compared to traditional TCP flavours.
- Suffers performance loss during timeout compared to N individual flows.
- More aggressive than should be for SACK, when $N < 5$.

High Speed TCP (HSTCP)

How?

It modifies the response function of TCP, so that ridiculously low drop rates are not required to sustain high throughput.

Increase function : $w = w + a(w) / w$

Decrease function : $w = w (1 - b(w))$

For $\text{cwnd} < 38$, $a(w) = 1 + b(w) = 1/2$ (same as standard TCP)

For larger cwnd, $a(w)$ & $b(w)$ are long functions based on several measures. $a(w)$ & $b(w)$ are precomputed and stored in lookup tables.

Observations:

HTCP emulates roughly N parallel TCP connections.

Scalable TCP (STCP)

(11)

Modifications to TCP Algorithms:

- For $cwnd \leq lwnd$, we traditional window update algorithm.
- For $cwnd > lwnd$,
 - Increase function = for each ack, $cwnd = cwnd + 0.01$
 - Decrease function = for each loss, $cwnd = cwnd - [0.125 * cwnd]$
- $lwnd$ chosen to be 16 corresponding to drop probability 5.25×10^{-3}
- This response function efficiently PS MIMO.
- Fixed time for doubling sending rate (TD RTTS).

Comparison with Traditional TCP:

<u>Traditional TCP scaling properties</u>	
For each ack $\rightarrow cwnd = cwnd + 1 / cwnd$	<u>Scalable TCP scaling properties</u>
For each loss $\rightarrow cwnd = cwnd - 0.5 * cwnd$	$cwnd = cwnd + a$

Scalable TCP scaling properties

$$cwnd = cwnd + a$$

$$cwnd = cwnd - [b * cwnd]$$

Why $a=0.01$ & $b=0.125$?

- Impact on legacy traffic
- bandwidth allocation properties
- flow rate variance
- Convergence properties
- control theoretic stability

Pros:

- A special case of HSTCP response function
 - must be easy to implement, model mathematically & less ad-hoc.
- Fixed doubling time for sending rate (scalable)
 - Independent of link capacity

Cons:

- MIMD inherently more aggressive than SIMD
- Not enough experimental results to show TCP will not starve standard TCP.

(12)

Fast AQM scalable TCP (FACT)

- Built on the TCP-Vegas core.
- Uses equation based control with queuing.
- Uses multiplicative decrease with packet loss.
- Based on flow-level dynamics rather than packet level dynamics.

Main Algorithm

- Estimate target rate
- Estimate difference between current rate & Target rate
- If very high, increase sending rate aggressively
- If not, increase sending rate very smoothly.
- Avoid oscillations, keep router queues stable.

Window Control Algorithm

For each RTT,

$$W = w(\text{base RTT}/RTT) \alpha$$

where,

RTT = exponential weighted average of RTT

base RTT = minimum instantaneous RTT

α = constant determining fairness / aggressiveness & convergence rate

α is computed based on link bandwidth, 1Pkt RTT etc.

If it is stored in a look up table.

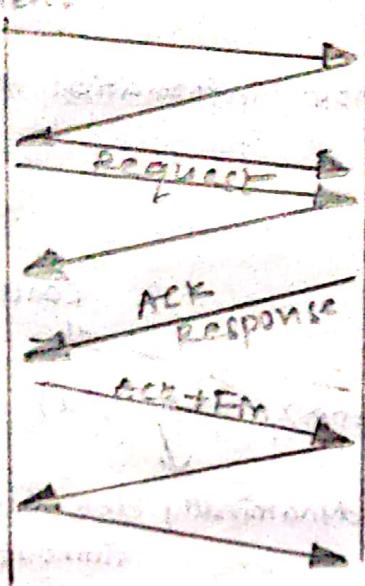
Pros & Cons

- ④ Theoretically fair, not enough experimental results.
- ④ Delay for congestion indication not fully understood or widely deployed yet.

T/TCP : Transaction Oriented TCP

(13)

Avoids the 3 way handshake to prevent long delays for transaction-oriented client server applications.



Why T/TCP?

- standard TCP uses atleast 10 segments for a single segment request - response sequence.
 - high relative overhead
 - Extra transaction latency for client on CPU & network speed increase.
- The 3 way handshake precludes high performance transaction processing
- A high rate of short connections like transactions leads to excessive consumption of kernel memory by connection control blocks in TIME WAIT state.

How does it work?

Transaction Oriented TCP

- suppresses the 3 way handshake
- shortens TIME WAIT state

Bypassing the 3-way handshake.

TCP Accelerated open (TAO) Mechanism:

- Validates initial SYN segments
- Enforces at most one semantics without a 3-way handshake.

Basis: TCP uses cached per host information to immediately validate new SYNS.

Shortening TIME-WAIT state

- Set TIME-WAIT delay to

$$\text{TIME-WAIT-Delay} = \max(k * RTO, u)$$

constants

dynamically determined retransmission timeout

- If RTO is very small, short connections are possible.

- TIME-WAIT delay reduces to u.

To accelerate close sequence, u has to reduce below MSL set by IP layer.

- Done using a monotonic number sequence, x.

- To further reduce TIME-WAIT-Delay

- explicitly truncate TIME-WAIT state

- remove TIME-WAIT Delay entirely.

TCP Performance Issues Over Wireless Links

TCP Flow Control

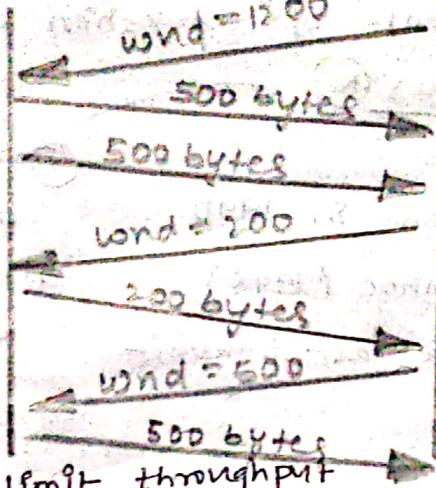
- prevent buffer overflow at receiver
- useful for fast sender transmitting to slow receiver.

How does it work?

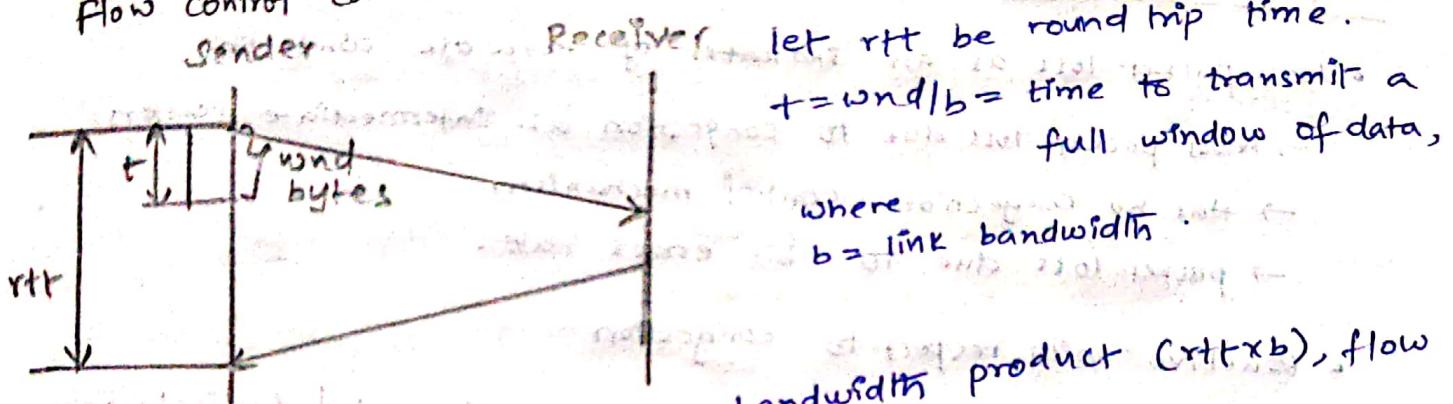
- Receiver advertises a window (wnd) in ack returned to sender.

- Sender cannot send more than wnd unacknowledged bytes to receiver.

(15)



flow control can limit throughput



For a link with high delay-bandwidth product ($rtt \times b$), flow control window can limit connection throughput

$$t \leq rtt, \text{ throughput} = \frac{wnd}{rtt}$$

TCP Congestion Avoidance

- Reduce congestion inside network: Requires cooperation of multiple senders, must rely on indirect measures of congestion.

- Implemented at sender

- Introduce a congestion window (cwnd) in addition to flow control window (wnd). size of this window needs to be managed.

Slow Start

- Aggressively grow congestion window until congestion is detected.

Loss Detection & Retransmission

- Fast retransmission & recovery
- less severe adjustment of congestion window size

Options

"TCP Reno"

- Developed by Van Jacobson in 1990
- Improvement to TCP Tahoe (1998)
- Added fast recovery & fast transmit

④ At end of options 184, max. segment

④ Later RFCs, are defined.

Congestion Avoidance

- Uses packet loss as an indicator of network congestion.

→ Most packet loss due to congestion at intermediate routers.

→ Has no congestion control mechanism

→ packet loss due to bit errors rate.

- Reactive with respect to congestion

→ Responds to packet losses indicated by timeouts/duplicate ACKs.

- Packet loss

→ Reduce congestion window size

→ Due to timeout: cwnd = 1 and ENTER SLOW START

→ Due to duplicate ACKs: $cwnd = \frac{cwnd}{2} + 3 * \text{segment_size}$.

- Congestion window size increased when data successfully acknowledged.

- slow start:

- Active if $cwnd \leq ssthresh$ (threshold)

- During slow start, (congestion window increased by segment size for every ACK received, it opens window exponentially)

- Congestion Avoidance:

- $cwnd = cwnd + (1/cwnd) + \text{segment_size}/8$ for every ACK received

- Additive growth in window size, about 1 segment every rtt.

Slow Start Mechanism

- Detects, avoids congestion as connection begins / after a timeout.
- RENDE doubles congestion window every round trip if no loss occurs
- Aggressively reduce rate when invoked

Loss Detection

- coarse grained timeout: indicates packet loss
 - sender starts timer when TCP segment sent
 - Timeout if ACK not received before timeout
 - Retransmission
 - slow start invoked
- 3 duplicate ACKs indicate packet loss
 - Receiver to send ACK out of older segment
 - sender to assume loss when 3 duplicate ACKs received
 - Fast transmission & recovery

Pros:

- Simple bandwidth estimation scheme
- Aggressive congestion avoidance mechanism
 - ↓
 - Ensures bandwidth when connected to TCP Vegas.
- More widely deployed

Cons:

- periodic oscillation in window size
- oscillation in RTTs
- inefficient bandwidth utilization
- Many retransmissions of same packets after a packet - dropped
- Biased against connections with longer delays
- Unfair share of network capacity

TCP Vegas

Developed by Brakmo & Peterson - 1995

New congestion avoidance algorithm

Congestion Avoidance

Set congestion window size \rightarrow difference b/w expected & actual data rates

\rightarrow To control no. of outstanding bytes in queues in network

$$\text{diff} = (\text{cwnd}/\text{rtt}^* - \text{cwnd}/\text{rtt}) \times \text{rtt}^*$$

Where,

cwnd = current congestion window size

rtt* = minimum rtt

rtt* = actual rtt

diff = estimate backlog in queue

α = low threshold for diff ($\text{diff} > \alpha$)

β = high threshold for diff ($\text{diff} < \beta$)

Vegas keeps atleast α bytes but $< \beta$ bytes in queue.

- If $\text{diff} < \alpha$, cwnd + 1
- If $\text{diff} > \beta$, cwnd - 1
- Else ($\alpha \leq \text{diff} \leq \beta$), cwnd not changed.

Vegas provides proactive response to congestion. Congestion window changed gradually as observed backlog changes.

Slow Start Mechanism

Doubles congestion window every other round trip if no loss occurs.

Loss Detection

- Coarse grain timeout mechanism
- Fine grain timeout mechanism
 - Duplicate ACK received & RTT of 1st unacknowledged segment > fine grain timeout value
 - Non duplicate ACK received after retransmission & RTT of segment > fine grain time out value.

Pros:

- Fair bandwidth estimation scheme
- smooth sending rate, efficient link utilization
- Detect losses faster than TCP Reno & recovers more efficiently from multiple drop

Cons:

- Cannot compete with more aggressive Reno connections
- May not be stable if buffers are small

TCP Reno Vs TCP Vegas

In a homogenous environment, Vegas outperforms Reno.

40% - 70% better throughput

20% - 50% of losses

Factors:

- slow start & congestion avoidance have greatest influence on throughput
- Congestion detection mechanism has only minor effect
- May exhibit problems with fairness among competing connections and conditions.

TCP problems in wireless Networks

- Packet loss due to
 - Bit errors due to wireless channel impairments
 - Handoffs due to mobility
 - Possibly congestion

- TCP assumes packet loss due to
 - Congestion in network
 - Pack reordering

- TCP congestion avoidance can be triggered by packet loss
 - Do not respond well to packet loss due to bit errors, handoffs
 - performance is based on application, can suffer.

- Bursts of error due to low signal strength

- More than 1 packet loss in TCP
 - More likely to be detected as timeout
 - enter SLOW-START

- Delay often very high

- RTT can be very long & variable
 - TCP's timeout mechanism may not work well
 - Made worse by link level retransmission

- Links may be asymmetric.

Solutions to these problems

Link-layer approaches

- Hide losses due to congestion from sender
- make link appear more reliable
- ARQ, forward error correction (FEC) code
- Hybrid ARQ & FEC

Pros

- No change to existing sender behaviour
- Matches layered protocol model

Cons

- Interaction with TCP,
Ex: fast transmission with TCP can be triggered by delays due to link layer timeout & retransmission.
- Negative interactions with TCP reduced by making link-level protocol TCP aware.

Ex: Snoop TCP

Advantages

- Attempts to retransmit & suppress duplicate ack.
- state is soft
- Hand off simplified

Disadvantages

- May not completely shield TCP from effects of mobility.
→ wireless links.

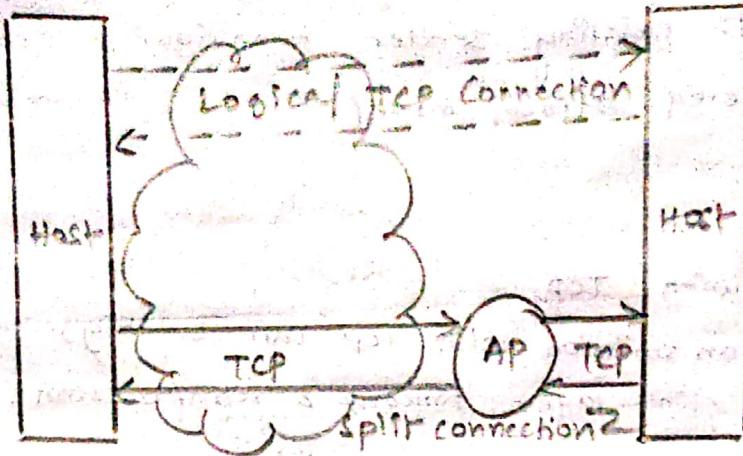
Split Connection Approaches

- Hide wireless link entirely by terminating TCP connection prior to wireless link.
- Hides at base station / access point
- uses special protocol / regular TCP over wireless link

Ex: Indirect TCP

Problems

- Extra protocol overhead
- violates end-to-end semantics of TCP
- complicates hand off
- state info. at AP / base where protocol is split



End-to-End approaches

- Make TCP sender aware of some issues not due to congestion.
- Avoid congestion control when not needed.
- Use selective acknowledgement (SACK) for "fine-grained" error recovering
Ex: SACK RFC, SMART
- Use explicit loss notification (ECN) to distinguish between congestion & other losses.

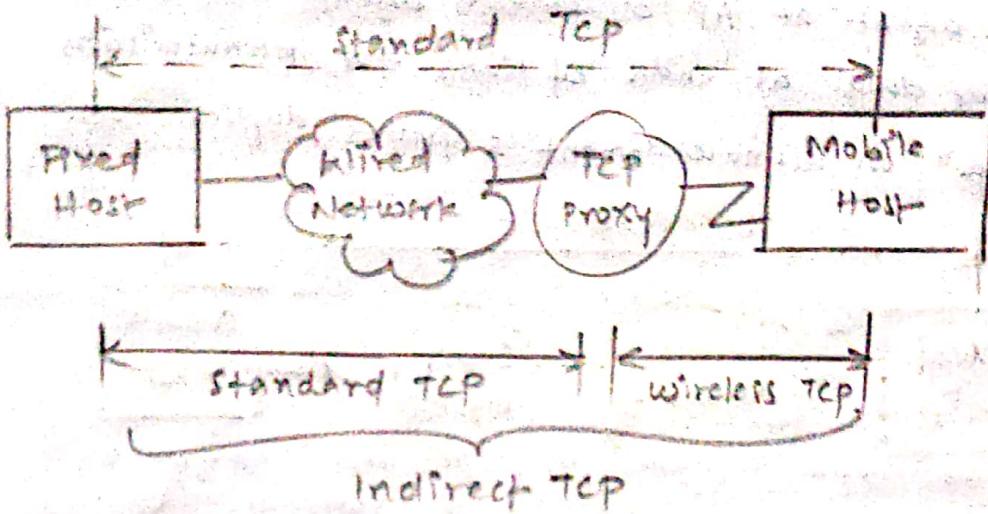
Advantages

- Maintains end-to-end semantics of TCP
- No extra overhead at base stations for protocol processing/hand off.

Disadvantages

- Requires modified TCP
- May not operate efficiently
Ex: packet reordering vs packet loss.

Indirect TCP



Hand off

- Access point/router can act as mobile IP foreign agent & TCP proxy for indirection TCP (I-TCP)
- Mobile host moves to different foreign agent (handoff needed for mobile IP)
- Mobile host moves to different proxy (handover/hand off of full TCP header for I-TCP)

Advantages

- No changes to TCP at hosts in fixed networks
- Errors from wireless links corrected at TCP proxy
- optimisation possible over wireless link

Disadvantages

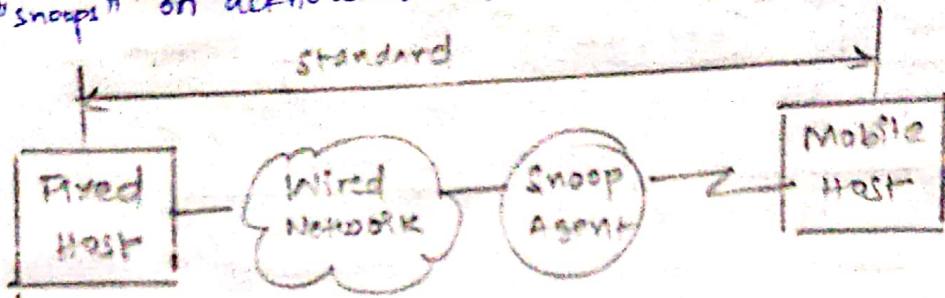
- loss of TCP end-to-end semantics
- Handoff overhead can be significant
- overhead at proxy for per-packet processing
- TCP proxy must be trusted

Wireless Transport

- uses TCP as wireless transport protocol
- Timeout at wireless sender may stall original sender
- selective acknowledgement protocols provide better performances

Snoop TCP

- provide reliable link layer, i.e. TCP-aware
- Snoop agent at AP or foreign agent
- Buffer data at ends of links for retransmission
- "Snoops" on acknowledgment & filters duplicate acknowledgment



Operation

- Snoop agent monitors & buffers data from fixed network to mobile host.
- Snoop agent monitors ACKs from mobile host.
 - can discard buffer data when acknowledged
 - can retransmit data (delayed/duplicate) ACK
- short timeout → fast retransmission
- Snoop Agent discards duplicate ACKs from mobile host
- Snoop Agent discards duplicate data already sent by agent & acknowledged
- Snoop Agent cannot generate ACKs sent back to fixed host
 - ↓
 - end-to-end TCP semantics preserved.

Reverse Direction

- Snoop monitors traffic from mobile host back to fixed host and detect missing segments.
- A negative ACK sent immediately to mobile host.
- Mobile hosts can retransmit missing segment, in time to avoid TCP time out at fixed host.

Advantages:

- Preserves end-to-end semantics
- No changes in TCP for fixed hosts
- No changes in TCP possible for mobile hosts; reserve direction traffic can benefit from days.
- No need for hand off
- Automatic fall-back to standard TCP

Disadvantages

- Does not fully isolate wireless link errors from fixed network
- Mobile host modified to handle NACKs for reverse traffic
- Cannot snoop encrypted datagrams
- Retransmission of data from agent not authenticated.

TCP Options

The only options defined in original TCP specification are the end of options list, no operation & maximum segment size option. Later RFCs 793 & 1323 defined new options.

→ End of option list

Kind = 0
1 byte

→ No operation

Kind = 1
1 byte

→ Max. segment size

Kind = 2	len = 4	MSS
1 byte	1 byte	2 bytes

→ Window scale factor

Kind = 3	len = 3	shift count
1 byte	1 byte	1 byte

→ Timestamp

Kind = 8	len = 10	Timestamp value (+ bytes)	Timestamp echo reply (+ bytes)
1 byte	1 byte		

*TCP proposal specifies 3 new options with kind = 11, 12 & 13.

④ They are connection count (cc), NEW, ECHO & cc.

⑤ Four other options of kind = 4, 5, 6 & 7. They are called 'echo-options'. They are replaced by Timestamp.

Stream Control Transmission Protocol (SCTP)

②

- IETF approved as a proposed standard in Oct 2020.
- A general purpose transport layer protocol
- Offers point-to-point, connection-oriented, reliable delivery transport service.
- Powerful congestion control & packet loss recovery
- Any app. running over TCP can be ported to run over SCTP
- Support for multihoming and partial ordering

enables SCTP host to establish "session" with another SCTP host over multiple interfaces, identified by separate IP addresses.

lets SCTP provide in-order delivery of one or more related sequences of messages following between 2 hosts

- Reliable delivery of fast processing of multiple, unrelated data streams.

TCP Issues

- Strict order of transmission delivery since, service.
- Treats each data transmission as an unstructured sequence of bytes
- No support for multihoming → socket-based application programming interface
- Susceptible to denial of service attack → TCP SYN storms.

SCTP features

Association → protocol state installed on 2 peer SCTP hosts exchange messages

- Message boundary preservation
- No head-of-line blocking

- Multiple delivery nodes
- Multihoming support
- TCP congestion control
- selective acknowledgments
- User data fragmentation
- Heartbeat keep alive mechanism
- DOS protection

Multihoming

- Offers network resilience to find failed interfaces on the host and fast recovery during network failures
- less affected by network reconvergences → lost packets retransmitted to alternate address
- Stress Reduction: Reduce stress on internet routing system
- Topology diversity: diversity in routing paths dictates level of fault tolerance
- Delivery options: separate reliable & ordered delivery into 2 independent functions.

SCTP Initiation

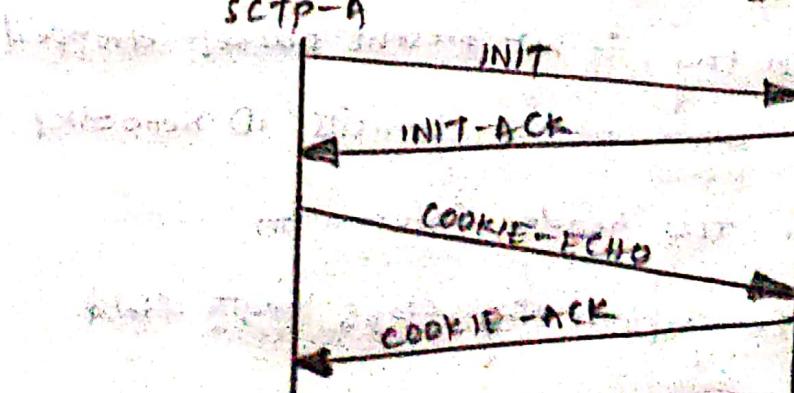
SCTP Association defined as

[a set of IP addresses at A] + [port A] + [a set of IP addresses at z] + [port -z]

Before data exchange, exchange communications state using a 4-way handshake. Eliminates exposure to TCP SYN flooding attacks

SCTP-A

SCTP-z



1. INIT (Initial Contact message)

2. INIT-ACK → with state cookie

Parameters in both:

- list of all IP addresses part of association
 - initial transport seq.no for reliable data transfer
 - initiation tag for every in bound SCTP packet
 - No. of outbound stream each side is requesting
 - No. of inbound streams each side can support
 - might have user DATA messages
3. COOKIE-ECHO → state cookie
4. COOKIE-ACK → sent after receiver fully reconstructs its state & acknowledges.
- setup is complete
 - can also bundle user DATA messages.

SCTP Data Transfer

1	Source port	Dst. port
2	verification Tag	
3	checksum	
	chunk-1	
	chunk-2	
	chunk-n	

Type	Flag	Length
TSN		Transport seq. no
Stream ID	SSN	Stream seq. no
Protocol ID		
User Data		

SCTP packet Format (common header + one/more variable length data chunks)

- 1 → Identify recipient of sctp packet
- 2 → value of initiation tag, if not present packet dropped
- 3 → Assure integrity while packet transits ID networks.

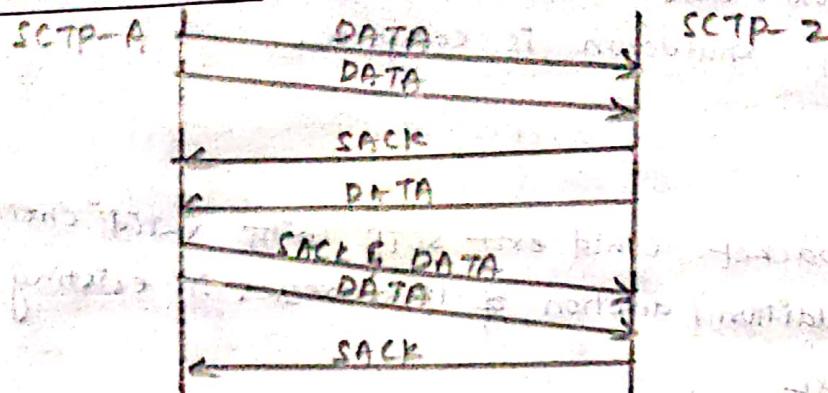
Every chunk includes TLU header information

- chunk type / delivery processing flags, length field.

TSN → for per-association reliability

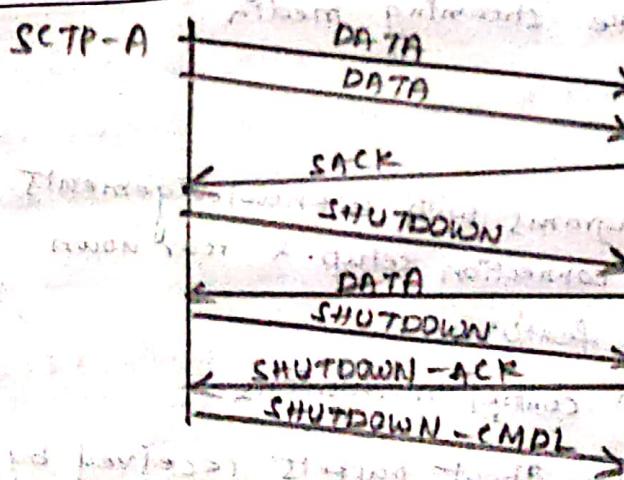
SSN → for per-stream ordering

SCTP Data Exchange



SCTP host sends SACKs in response to every SCTP packet carrying data chunk. The SACK fully describes receiver's data.

SCTP Shutdown



No "half-closed" state in SCTP

A wishes to shutdown & terminate association with 2.

1. SCTP enters SHUTDOWN-PENDING STATE

→ accept no data from app, send no queued data to 2.

2. After acknowledging all queued data, A sends SHUTDOWN chunk & enters SHUTDOWN-SENT state.

3. upon receiving SHUTDOWN chunk, 2 notifies upper layer

→ stops accepting new data

→ enters SHUTDOWN-RECEIVED state.

4. 2 transmits any remaining data to A
5. SHUTDOWN chunks follow to inform 2 → data arrived @
assoc. shutting down
6. All queued data on 2 acknowledged → A sends
SHUTDOWN-ACK chunk, then, SHUTDOWN-COMPLETE chunk.
7. Association shutdown is complete.

Future Issues

- Corrupted packet could exit SCTP with valid checksum
- Dynamic addition / deletion of IP addresses in existing association
- DPRT support

Datagram Congestion Control Protocol (DCCP)

- Congestion controlled, unreliable flow of datagrams
- Suitable for apps like streaming media

DCCP Features

- Unreliable flow of datagrams with acknowledgements
- Reliable handshake for connection setup & tear down
- Reliable negotiation of features
- TCP friendly congestion control mechanisms
- Options to inform sender about packets received by receiver & their status
- Congestion control incorporating explicit Congestion Notification (ECN)
- Allow server to avoid holding state for unacknowledged / already finished connection.
- Path MTU discovery.

Important Differences from Tcp

- Packet stream protocol, not byte stream
- Unreliable → no transmission of datagram

- packer seq. nos, not for bytes
- space for options (upto 1020 bytes)
- Feature negotiation
- choice of congestion control (2 end-points)
- Different acknowledgment format
- No receive window → congestion control, not flow control protocol
- Distinguishing different kinds of loss
- Definition of acknowledgement (which options have been processed)
- Integrated support for mobility
- No stimulations open

What is a Half-connection?

- A subset of connection
- Data packets sent in one direction + corresponding ack. sent in other direction
- HC sender : sends data
- HC receiver : sends ACKs.
- Every connection (half) managed by congestion control mechanism.
- specified by CCID congestion control identify.

Typical PCCP Connection

1. DCCP Request: client → server, service, features, CCID req.
2. DCCP Response: server → client, indicates willingness to communicate
3. PCCP-ACK: client → server, acknowledges the response packet

4. Of more DCEP-Ack : To finalise feature negotiation
5. Exchange of DCEP-Data, DCEP-Ack and/or DCEP-Data ACK.
6. DCEP-Close Req : server sends request to close.
7. DCEP-Close : client acknowledges the close.
8. DCEP-Reset : server sends reason set to "closed" & (server) cleans connection state.
9. DCEP-Reset : Client receives packet & holds state for any (client) remaining packets to clear.