

lecture 5/

Transport layer

Overview

- Application layer (web, mail, etc)
- • Transport layer (process-to-process data transfer)
 - Network layer (transports pkts between nodes in network) IP, etc
 - Link layer (transports pkts between neighbouring nodes)
encoding, framing, shared medium, reliable data transp.
 - Physical layer (transports bits between neighbouring nodes)
info theory, physical media, Shannon capacity

Types of end-to-end services and protocols

- unreliable pkt datagram delivery → UDP
- reliable bytestream service → TCP
- request/reply service → RPC (Remote Procedure Call)
- real-time service (no guarantee of timely delivery) →
→ RTP (Real-time Transport Protocol)

UDP

- Trivial host-to-host protocol (on top of IP)
- ↳ Each application pkt is sent as an independent IP pkt
no connection setup
- IP's best effort service: pkts may be lost/reordered
- Packet size limited to IP capacity
- Provides multiplexing/demultiplexing
- Provides error detection using simple checksum
- message format:

port (

- src 16
- length 16
- data 32
- dst 16
- checksum 16

- connectionless (demultiplexing)

- TCP
- Provides reliable, in-order delivery (using ACKs & retransmissions)
- Provides multiplexing (demultiplexing)
- Full-duplex
- Connection-oriented (different sockets per each connection)
- Flow-control: prevents sender from overrunning receiver
- Congestion-control: prevents sender from overloading the network
- Byte-oriented: apps write/receive stream of bytes (although network transports them in packets framing in app. layer)

- End-to-end issues
- At link layer, we used retransmissions to make it reliable
- Differences from end-to-end protocols:
 - RTT is not known beforehand (varies and changes)
 - pkts may get reordered in the Internet (TTL=120s)
 - bandwidth-delay product and buffer space at peer are not known beforehand
 - max. packet size on entire path is unknown

• TCP Header

src port 16 | seq. num (first byte in segment) 32 | ~~data~~

dst port 16 | seq. num (of) |

• TCP Header

src port 16 | seq. num 32 | header length 4 |

dst. port 16 | Ack num 32 | unused 6

• flags:

URG 1 | seq. num of next byte

ACK 1 | advertised window size 16 | options (var lengths) 3

PSH 1 | TCP checksum 16 | data

RST 1 | urgent pointer 16

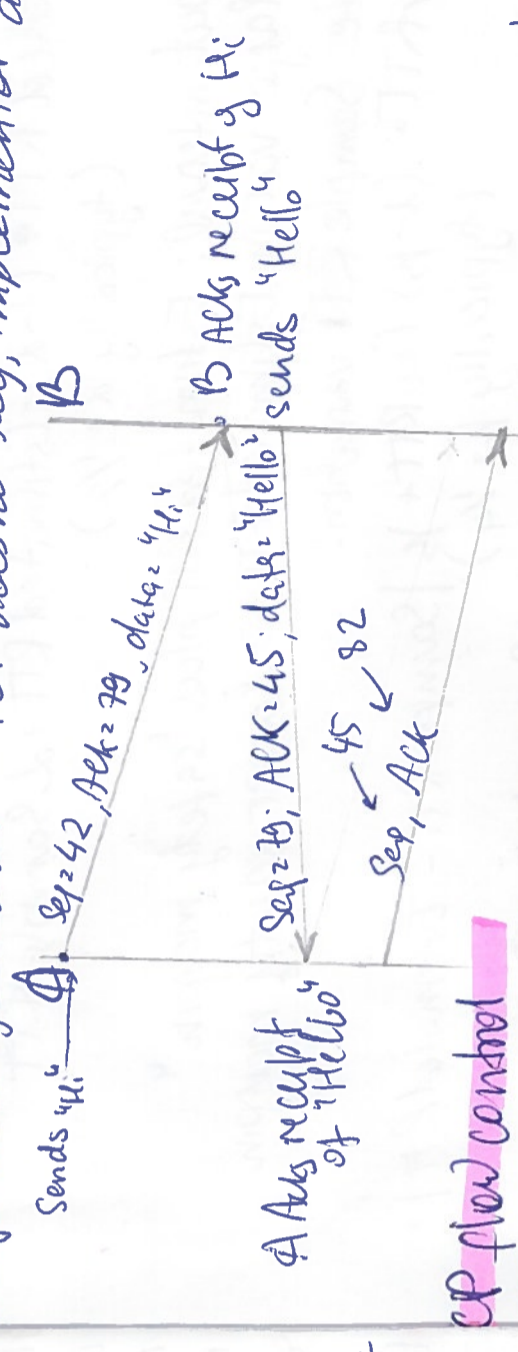
SYN 1

FIN 1

TCP sequence numbers and ACKs

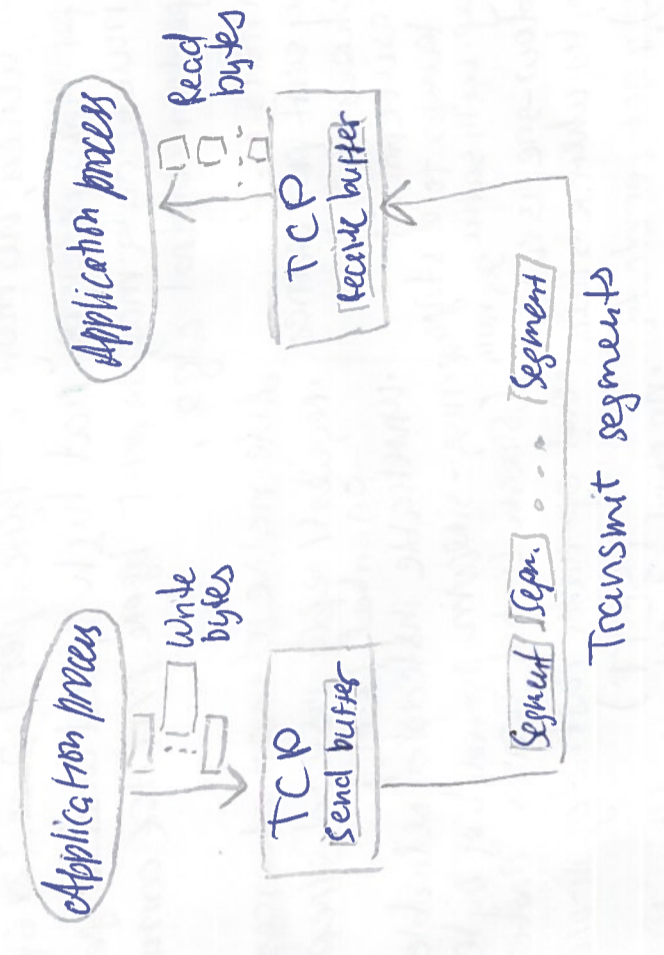
- sequence numbers:
- byte stream "number" of first byte in segment's data
- ACKs
- Seq. nr. of next byte expected from the other side
- ACKs are cumulative

out-of-order segments? — TCP does not say, implementor decides



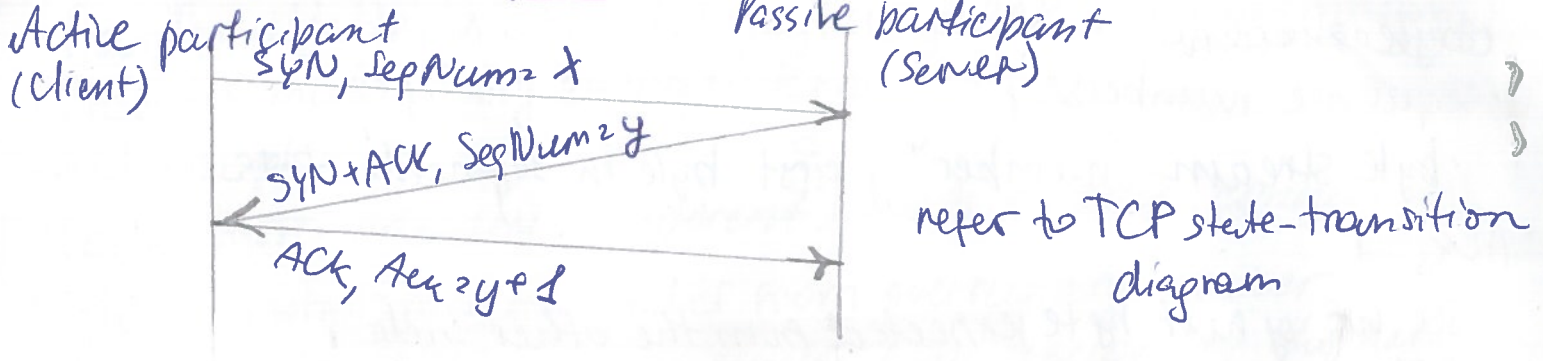
TCP flow control

- Goal: prevent sender from sending faster than receiver can handle
- Principle:
 - Receiving TCP advertises remaining buffer space ("advertised window" in TCP Header)
 - Sender is not allowed to send more data than fits in receiver's window



- If receiving app stops
- Receiver:
 - receive buffer fills up
 - advertised window decreases, to zero
 - sending TCP cannot send data
 - send buffer fills up
 - Sending application will be blocked if no buffer space left
 - backpressure

TCP connection setup



refer to TCP state-transition diagram

"three-way handshake"

When to transmit?

Estimated RTT = $(1-d) \cdot \text{Estimated RTT} + d \cdot \text{Sample RTT}$
 (typically $d = 1/8$)

timeout interval: Estimated RTT plus "safety margin"

• large var in Estimated RTT \rightarrow larger safety margin

estimate Sample RTT variation:

Dev RTT = $(1-\beta) \text{Dev RTT} + \beta \cdot |\text{Sample RTT} - \text{Estimated RTT}|$
 (typically $\beta = 1/4$)

• then: Timeout Interval = Estimated RTT + $4 \cdot \text{Dev RTT}$

• do not use retransmitted pkt for RTT measurements

Problems with TCP at high speed

• Sequence number wrap: TCP/IP ttl is 120s. Seq num has 32 bits. To guarantee uniqueness, no more 2^{32} bytes per 120sec. (286 Mbps)

• window too small for "long fat pipe": at high speed and large transport delay much data may be sent before first ACK comes in

• problems with congestion control algo

TCP Header extensions

• Timestamp: every sent pkt carries timestamp, ACK echoes this; accurate RTT measurement

• PAWS: same timestamps to distinguish old from new seq. with same seq num

• window scaling: window size is not in bytes, but as 2^k bytes, where k is num.

• Selective ACK (SACK): ack for data that has been received out of sequence.

Alternative design choices

• request/reply instead stream-oriented (RPC)

• unreliable instead of reliable (RTP) msg-stream instead of byte-stream (Stream Control Trans. Protocols)

• out-of-order instead of in-order delivery (SCTP)

• no explicit set-up/teardown phases
 • rate-based instead of window-based