# Systems and Devices 2 (Network)
# Lec 3b: Transport Layer

---

# Before we get started ...

- Did you manage to teleport Bob (file transfer)?
  - Could you get my most excellent code working :)
- How did your research into TCP go?
  - Implementing a reliable communications protocol is tricky, need to consider all possible scenarios: lost / corrupt packets, packets received out of order, re-transmission of packets, fair usage of network bandwidth i.e. congestion control …
- Questions to consider :
  - How does TCP implement a reliable connection?
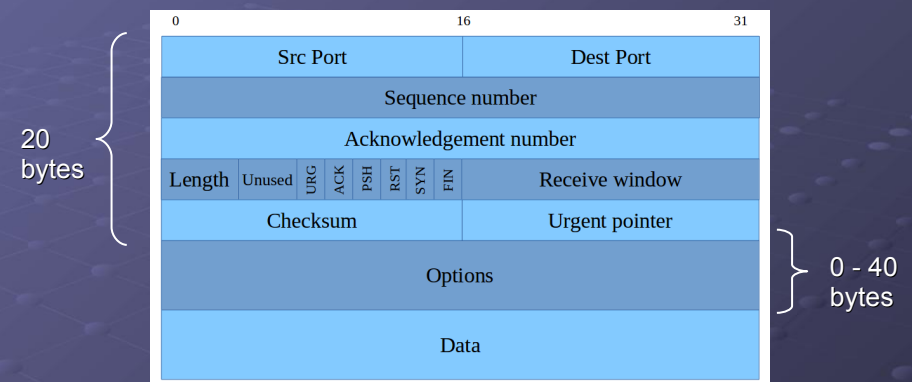  - What are the differences between UDP and TCP?

---

# TCP

- Transmission Control Protocol
  - RFC 793 : https://tools.ietf.org/html/rfc793
  - Created in 1981 one of the core internet protocols.
  - TCP is a connection orientated protocol i.e. before an application process can send data it must first perform a handshake to ensure that a connection is possible
    - TCP is <u>not</u> an end-to-end protocol like circuit switching, it only existing in the transport layer of the communicating hosts, not in lower layers i.e. no reserved connections, "packets" can take different paths.
  - TCP uses full duplex connections. If there is a connection between process A and B, then we can send data A→B and B→A at the "same" time.
  - TCP implements a reliable, ordered, error-checked data stream between two hosts i.e. retransmission.

---

# TCP



- TCP header, depending on option fields used can vary from 20-60B (UDP uses 8B).

# Demo



- TCP server: convert lower case string to upper case.
  - An echo function : **single**, multiple and continuous.

# Three-way handshake



- To initialise a TCP connection the two communicating processes perform a three-way handshake.
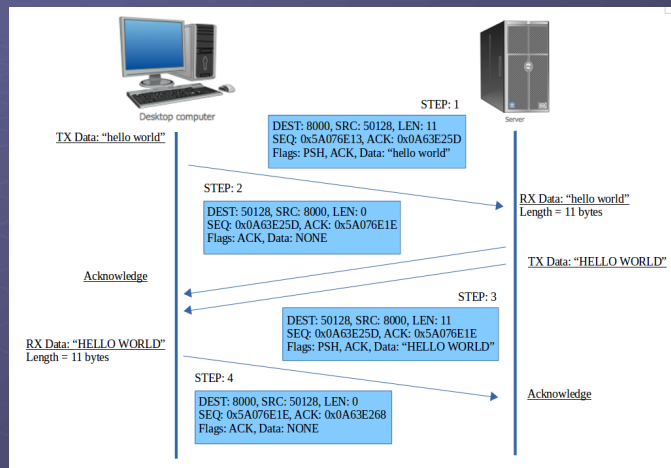
# TCP



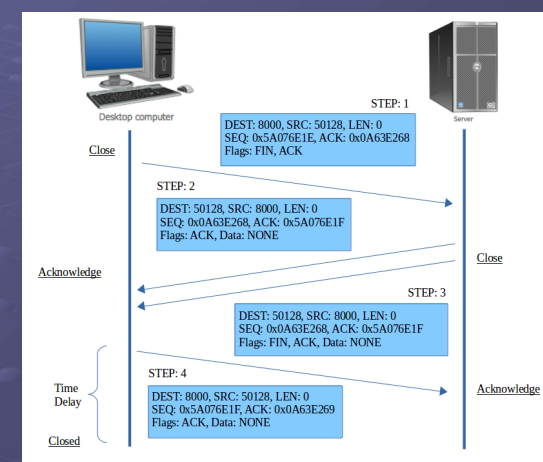- To transfer data the client uses the PSH flag

# TCP



- To close a TCP connection client/server use the FIN flag

## Pause to consider ...

- To ensure reliable connections TCP adds additional initialisation and acknowledgement segments. However, these incur significant overheads i.e. time.
  - ▶ To transfer the 10 characters UDP needs 1 segment, TCP needs 9 segments. Therefore, when using TCP we normally keep the link open and transfer multiple values across it.
- This highlights one of the main differences between TCP (stream) vs UDP (message) from the programming point of view i.e. how to separate out different values.
  - ▶ UDP: send individual messages i.e. paired send / receive
  - ▶ TCP: programmers responsibility to structure the data i.e. either by using a fixed length packets or delimiters such as newline characters.
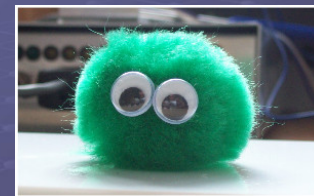
## Points to note ...

- Reliable data transfer is implemented using positive acknowledgements and timers (discussed next).
- Establishing a TCP link takes time i.e. 3-way handshake. If the client / servers are situated on different continents this will result in significant delays.
- Therefore, web browsers will tend to keep open TCP connection to avoid having to re-establish a link e.g. HTTP communicates across TCP links, using:
  - ▶ Non-Persistent : closed after each object is transferred.
  - ▶ Persistent : multiple objects sent over the same link.
  - ▶ HTTP1.1 : all connections assumed persistent, however, a web server's default time-out can be 5 – 15 sec.

## Demo



```
1 import socket
2 import time
3
4 TCP_IP = "192.168.0.254"
5 TCP_PORT = 8000
6 #TCP_PORT = 80
7 BUF_SIZE = 1024
8
9 sockTX = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
10 sockTX.connect((TCP_IP, TCP_PORT))
11
12 try:
13    while True:
14        lowerCase = "hello world"
15        sockTX.send(lowerCase)
16        upperCase = sockTX.recv(BUF_SIZE)
17        print("TX:" + upperCase)
18        time.sleep(2)
19 except KeyboardInterrupt:
20    sockTX.close()
21
22
23
```

- TCP server: convert lower case string to upper case.
  - ▶ An echo function : single, **multiple** and **continuous**.

## Programming Task



```
P3
# CREATOR: GIMP PNM Filter Version 1.1
80 50
255
248
247
246
224
215
```

- Q: reimplement the ppm image file (text file) transfer problem (code on VLE) using TCP, that should be simple :).
  - ▶ You will find this tricky. You need to think in terms of a stream of data, rather than packets of data.
  - ▶ Consider how the name, number of segments and image data will be sent across the TCP connection. Do you need to send the number of segments?

## Quick Quizzz

| | | |
|---|---|---|
| DEST: 8000, SRC: 50128, LEN: 10<br>SEQ: 0x5A076E13, ACK: 0x0A63E25D<br>Flags: PSH, ACK, Data: "abcdefghij" | DEST: 8000, SRC: 50128, LEN: 10<br>SEQ: _____, ACK: _____<br>Flags: PSH, ACK, Data: "abcdefghij" | DEST: 8000, SRC: 50128, LEN: 10<br>SEQ: _____, ACK: _____<br>Flags: PSH, ACK, Data: "abcdefghij" |
| SEGMENT 0 | SEGMENT 1 | SEGMENT 2 |

| | |
|---|---|
| DEST: 50128, SRC: 8000, LEN: 0<br>SEQ: _____, ACK: _____<br>Flags: ACK, Data: NONE | DEST: 50128, SRC: 8000, LEN: 0<br>SEQ: _____, ACK: _____<br>Flags: ACK, Data: NONE |
| ACK 0 | ACK 1 |

- The most important parts of the TCP header are the SEQ and ACK numbers i.e. to ensure reliability.
  - SEQ number is the byte stream number of the first byte in the segment.
  - ACK number is the sequence number of the next byte the RX host is expecting from the TX host.
- Q : what are the missing the SEQ and ACK numbers if:
  - RX in order segment 0,1,2 etc, or RX out of order 0,2,1?

## Pause to consider ...



- That was the easy bit :), now need to consider the what ifs :(
- IpoAC: multiple packets in flight at the same time, these may be received out of order, packet lose may also occur ...

## TCP

- To help understand how TCP operates consider the following simplified TX pseudo code.
- Three main events
  - Application data
  - Time-out
  - Acknowledgement

```
NextSeqNum ← RAN()
PrevSeqNum ← NextSeqNum
segment ← []

while TRUE do
    case EVENT of
        application_layer_data:
            segment.create( data, NextSeqNum )
            if timer.state = stop then
                timer.state ← start
            end if
            networkLayer( segment )
            NextSeqNum ← NextSeqNum + len( data )

        time_out:
            networkLayer( segment.select( notACK, smallestSEQ ) )
            timer.state ← start

        rx_ACK:
            if ACK > PrevSeqNum then
                segment.remove( ACK )
                PrevSeqNum ← ACK
                if segment.NotACK() then
                    timer.state ← start
                end if
            end if
    end case
end while
```

## Time-outs

$$\alpha = 0.125$$
$$estRTT = (1 - \alpha) \times estRTT + \alpha \times curRTT$$
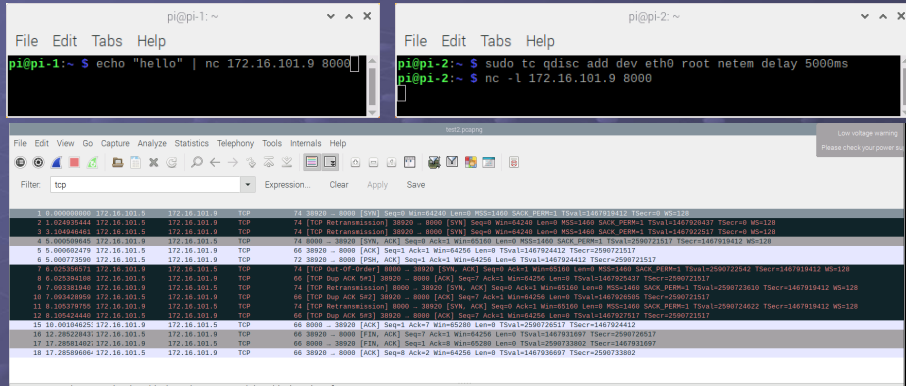
$$\beta = 0.25$$
$$devRTT = (1 - \beta) \times devRTT + \beta \times |\ curRTT - estRTT\ |$$

$$Timeout = estRTT + 4 \times devRTT$$

- Q: how do we choose the time-out value
  - Too short and we will reject valid segments in transit. Too long and we waste time waiting for segments that will never arrive.
- RFC 6298 : round trip time (RTT) delay. Initial time-out set to 1 sec. If time-out occurs, time-out value doubled, reset on next RX segment, do not use times from re-transmitted segments.

## Demo


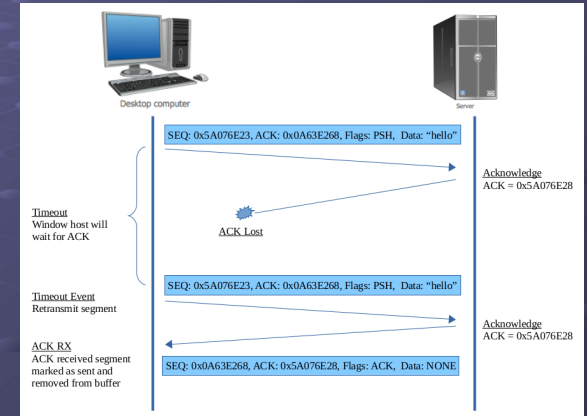
- Command line tool : tc (traffic control)
  - Add 5000ms delay from Pi-1 to Pi-2
  - Work through the time line of errors (black lines), what causes what, tricky :)
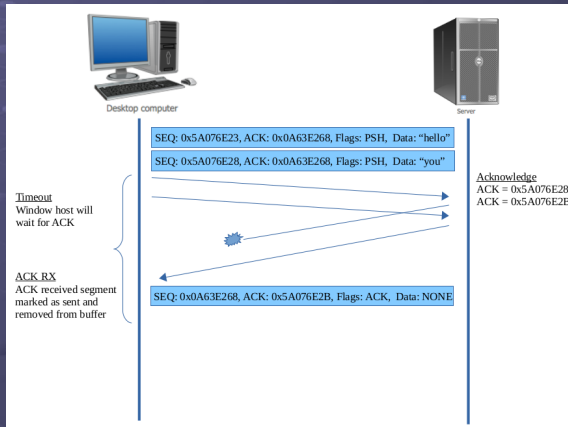
## Pause to consider ...

- Time-out example 1
  - TX host sends a message, it does not receive an ACK, timeout triggered.
  - Q: what happens next i.e. to the second "hello" RX by the server?
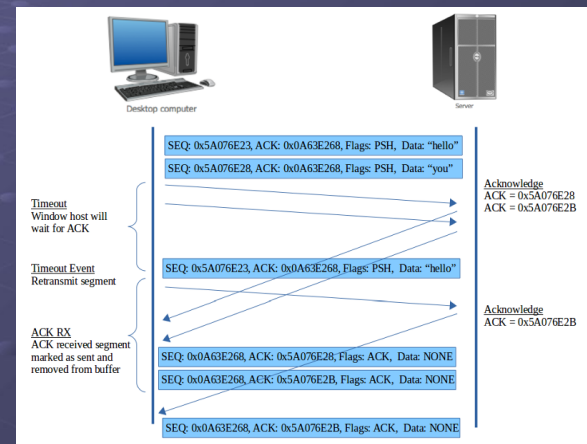
## Pause to consider ...

- Time-out example 2
  - TX host sends a message, but only one ACK received.
  - Q: what happens next i.e. does the client need the missing ACK?

## Pause to consider ...

- Time-out example 3
  - TX host sends a message, but ACKs delayed, triggering a timeout.
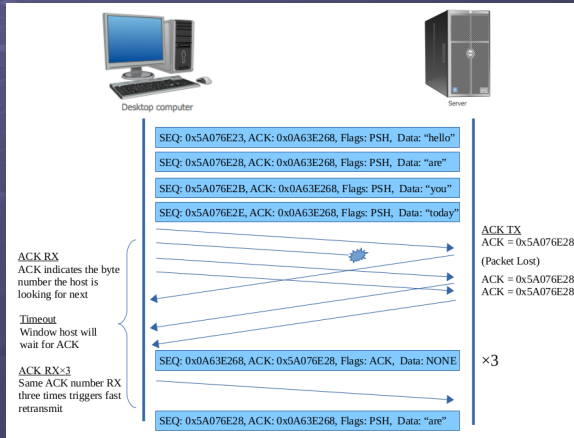  - Q: what happens next i.e. does the client know that all data has been RX?

## Pause to consider ...

- Time-out example 4
  - ► TX host sends a message but one of the "words" is lost.
  - ► Q: what happens next?
    - ◆ Fast retransmit
  - ► Q: what happens if only one word is sent?



```
Desktop computer                          Server

SEQ: 0x5A076E23, ACK: 0x0A63E268, Flags: PSH,  Data: "hello"
SEQ: 0x5A076E28, ACK: 0x0A63E268, Flags: PSH,  Data: "are"
SEQ: 0x5A076E2B, ACK: 0x0A63E268, Flags: PSH,  Data: "you"
SEQ: 0x5A076E2E, ACK: 0x0A63E268, Flags: PSH,  Data: "today"
                                                              ACK TX
                                                              ACK = 0x5A076E28
ACK RX                                                        (Packet Lost)
ACK indicates the byte
number the host is                                            ACK = 0x5A076E28
looking for next                                              ACK = 0x5A076E28

Timeout
Window host will
wait for ACK
ACK RX×3
Same ACK number RX
three times triggers fast
retransmit
SEQ: 0x0A63E268, ACK: 0x5A076E28, Flags: ACK,  Data: NONE     ×3

SEQ: 0x5A076E28, ACK: 0x0A63E268, Flags: PSH,  Data: "are"
```

---

## Pause to consider ...

- TX host transmits segments onto transport mechanism. Therefore, there can be multiple segments in flight.
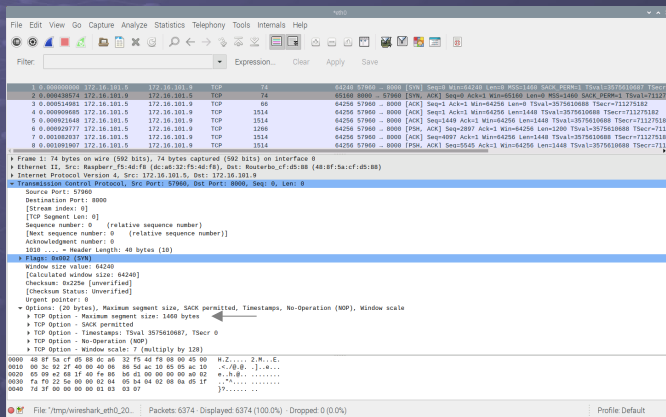  - ► $D_{End-End} = N_{stages}(D_{proc} + D_{Que} + D_{Tran} + D_{prop})$
  - ► Q : if we are transmitting 1000 byte segments to Australia at 1Mbps how many segments can be in flight?
- The RX host receives these segments, but may not process them immediately owing to:
  - ► Time taken to pass up the layers in the protocol stack
  - ► The OS is executing other processes.
  - ► There is a significant processing overhead for the application using these segments.
- Problem : RX host can be swamped with data i.e. data is coming in faster than it can process :(.
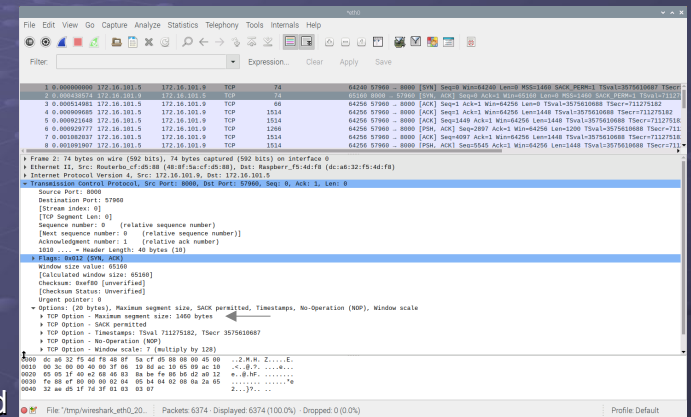
---

## MSS

- During 3-way handshake hosts signal MSS, max encapsulated data size
  - ► not including IP or TCP headers
  - ► If not specified defaults to 536B

---

## MSS

- During 3-way handshake hosts signal MSS, max encapsulated data size
  - ► not including IP or TCP headers
  - ► If not specified defaults to 536B

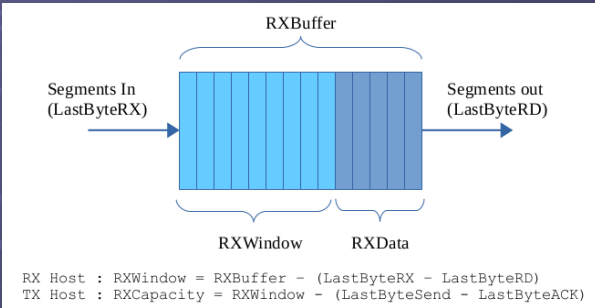# Flow control



Typical RXBuffer size set to Bandwidth Delay Product (BDP)

Typical LAN

1 Gbit/s * 1ms RTT
$10^9 * 10^{-3}$ = 125kB

Needs to be larger for slower networks
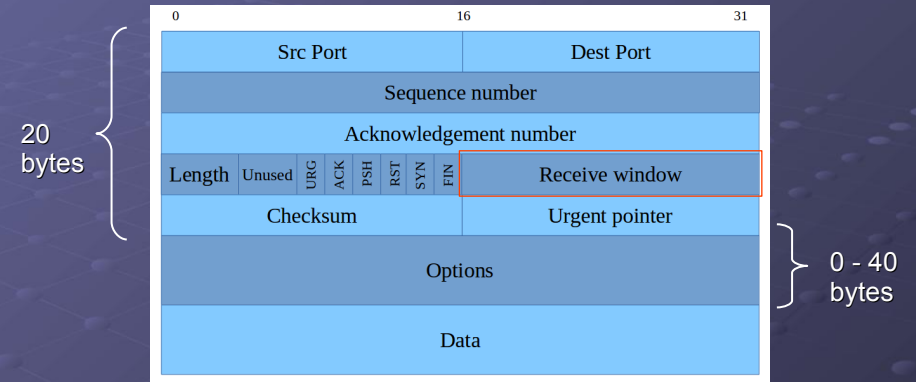
```
RX Host : RXWindow = RXBuffer - (LastByteRX - LastByteRD)
TX Host : RXCapacity = RXWindow - (LastByteSend - LastByteACK)
```

- Both TX and RX hosts have buffers to store data when busy. The free space in the RX host's buffer (RXWindow) is transmitted to TX host to ensure buffer overflow does not occur.
  - ▶ Data is temporarily stored on "wire" i.e. data in flight.
  - ▶ RX window size transmitted back to host in ACK packets.

---

# Flow control



- Hosts inform each other how much data they can receive using the Receive Window field (16bits)

---

# Flow control



- Window scaling factor sent during 3-way handshake

- Allowing RX windows greater than 64KB

---

# Flow control



- Receive window size : how much data can be in flight.

- Receive buffer size: how much data can be stored in RX host.

# Demo


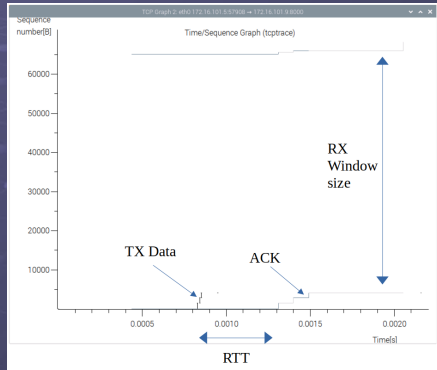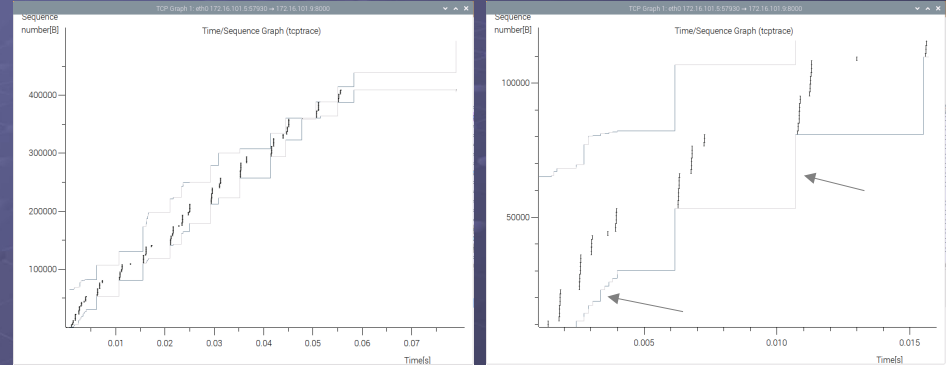
- Flow control : a little data (1).
  - ▶ Test code: simple python program transferring one or more TCP packets contain random values.

# Demo



- Flow control : a bit "more" data.
  - ▶ More data in flight
  - ▶ Move from individual ACKs to cumulative ACK

# Demo



RX buffer full
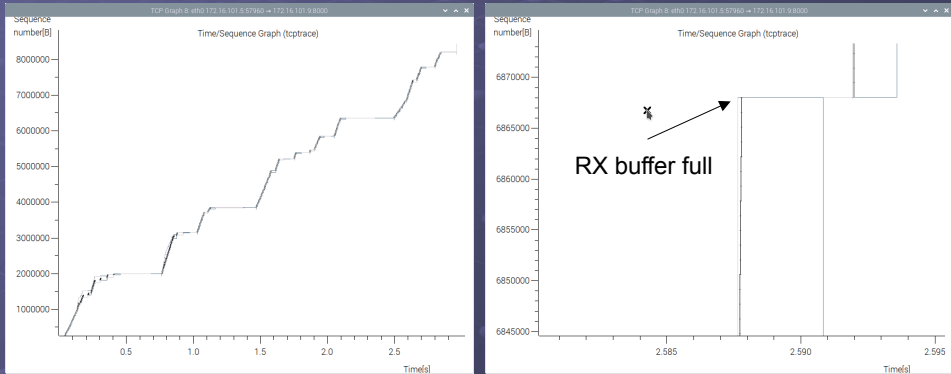
- Flow control : a lot of data.

# Demo



- RX window full (0) TX host must stop sending data.

# Pause to consider ...



- We have seen that TCP can adjust TX speed to match the RX speed to avoid buffer overflow.
- However, what should the system do if we just have too much traffic on the network i.e. segments are being lost because of buffer overflows in routers.
  - ▶ Remember packet switching uses store and forward.

---

# Congestion control

```
Congestion Window (CGWindow) = number of unACKed bytes

  Number of Bytes in transit    ≤ min( CGWindow, RXWindow )
( LastByteSent - LastByteACK ) ≤ min( CGWindow, RXWindow )

Assuming infinite RX buffer, Speed =  CGWindow ÷ RTT     (bytes per sec)

Initial Speed = MSS ÷ RTT      (then increase until missed ACK)
```

- TCP handles congested networks by using another "window", the Congestion window :). It detects a networks "loading" via the ACK segments
  - ▶ ACK received for unACKed segment. All is good increase congestion window size i.e. number of segment in flight.
  - ▶ Time-out, or multiple ACKs for same segment RX. Network congested, decrease congestion window size, i.e. packets are being dropped, reduce traffic.

---

# Congestion control

- TCP congestion-control algorithm
  - ▶ RFC 5681: https://tools.ietf.org/html/rfc5681
  - ▶ Three core elements: slow start, congestion avoidance and fast recovery (implementations do vary)
- Slow start :
  - ▶ Initially CGWnd set to 1 MSS i.e. speed= MSS/RTT.
  - ▶ On each RX ACK add 1 MSS to CGWnd i.e. doubled, speed increases exponentially, not too slow :).
  - ▶ Continues until ACK time-out i.e. host probes the network to find "max" speed. SSThresold set to CGWnd/2 and CGWnd reset to 1 MSS.
  - ▶ Process repeated until SSThresold reached, then switch to congestion avoidance mode.

---

# Congestion control

  - ▶ If three repeated ACKs for the same SEQ number RX, perform fast retransmit and change mode to fast recovery.
- Congestion avoidance:
  - ▶ CGWnd is half the value at which congestion was detected. Now rather than doubling CGWnd size it is incremented by 1 MSS when all CGWnd segments have been ACKed i.e. linear increase rather than exponential.
  - ▶ If ACK time-out. SSThresold set to CGWnd/2, CGWnd reset to 1 MSS, switch back to soft start mode.
  - ▶ If three repeated ACKs for the same SEQ number RX, perform fast retransmit and change mode to fast recovery.
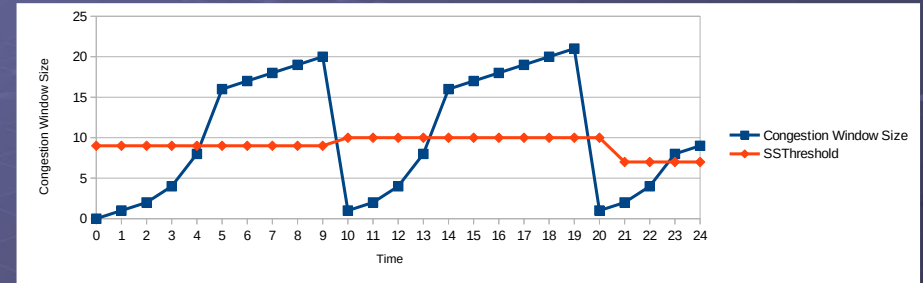
# Congestion control

- Fast recovery :
  - CGWnd set to (CGWnd/2)+3MSS. CGWnd increased by MSS for each duplicate ACK RX.
  - If ACK time-out occurs, CGWnd reset to 1 MSS and SSThresold set to CGWnd/2, switch back to soft start mode.
  - If ACK received, CGWnd set to SSThresold, then enter congestion avoidance mode.
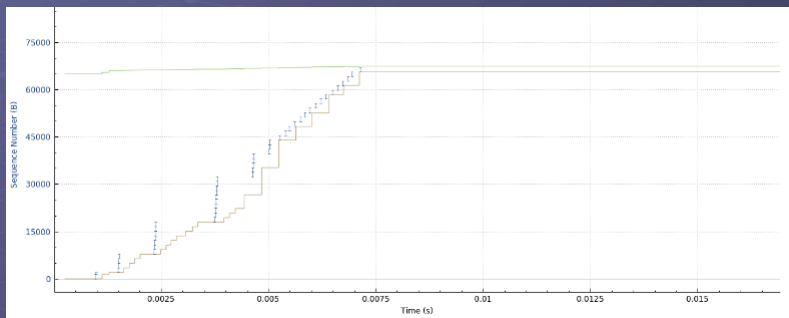
# Congestion control



- Note, different strategies taken depending on the version of TCP used. Refer to RFC5681 for additional details e.g. initial CGWnd size, strategies used in each mode to adjust CGWnd size etc.

# Congestion control



- Remember that flow control i.e. RXWindow size, is also active at this time, therefore, speed will also be throttled based on RX host buffer size.

# Summary

- TCP is used to implement:
  - Telnet, SSH, FTP, HTTP, SMTP, POP …
  - The internet's workhorse protocol.
- But, to ensure a reliable connection we need additional state information (memory on host), header fields and handshakes in protocol etc, i.e. reliability has a cost.
  - Sequence and acknowledgement numbers to ensure segments are not lost.
  - Time-outs, different types of segments (flags) i.e. SYN, ACK, PSH etc.
  - Initial three way handshake to ensure server is ready to communicate and agree upon parameters used in this transfer e.g. RXWindow sizes.

# Summary

- When you look into all the if-buts-and-maybes of how TCP works i.e. all the possible scenarios of time-outs, ACKs etc, its VERY complex.
- We have also not looked at the different option field in the header, or optimisations etc. If you have time you may want to have a read around the area.
- However, we still have some unanswered questions:
  - How does a host know how to route a TCP packet across the Internet i.e. we are not using circuit switching, so how do we know were to send our packets, how do we get from host A to host B?