



OST

Eastern Switzerland
University of Applied Sciences

Distributed Systems (DSy)

Protocols

Thomas Bocek

24.04.2024

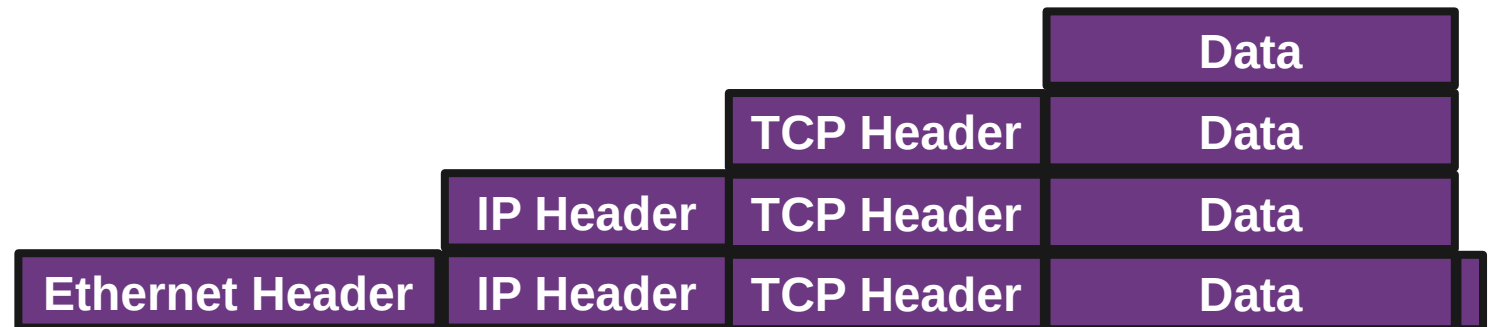
Learning Goals

- Lecture 8 (Protocols)
 - How do network layers work?
 - What are the TCP mechanisms?
 - What are the problems of TCP, and how other protocols (HTTP/3) can improve that

Networking: Layers

- Networking: Each vendor had its own proprietary solution - not compatible with another solution
 - [IPX/SPX](#) – 1983, [AppleTalk](#) 1985, [DECnet](#) 1975, [XNS](#) 1977
- Nowadays most vendors build compatible networks hardware/software from different vendors
 - Cisco, Dell, HP, Huawei, Juniper, Lenovo, Linksys, Netgear, MicroTik, Siemens, Ubiquiti, etc.
- Goal of layers: interoperability
 - 1984: ISO 7498 - The Basic Reference Model for Open Systems Interconnection

OSI model	"Internet model"
Application	Application
Presentation	
Session	
Transport	Transport
Network	Internet
Data link	Link
Physical	

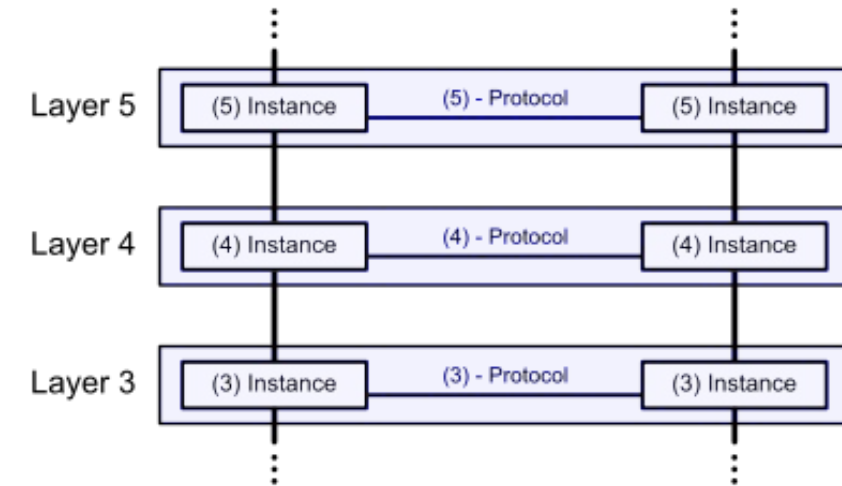
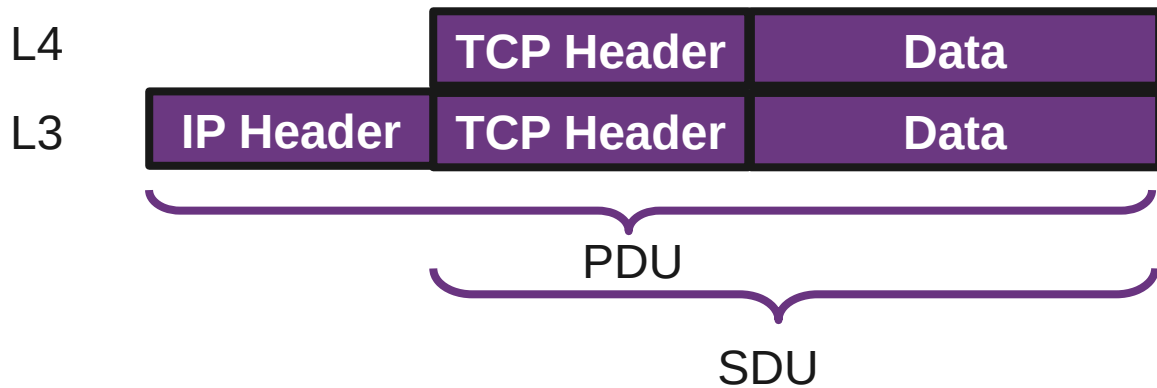


Networking: Definitions

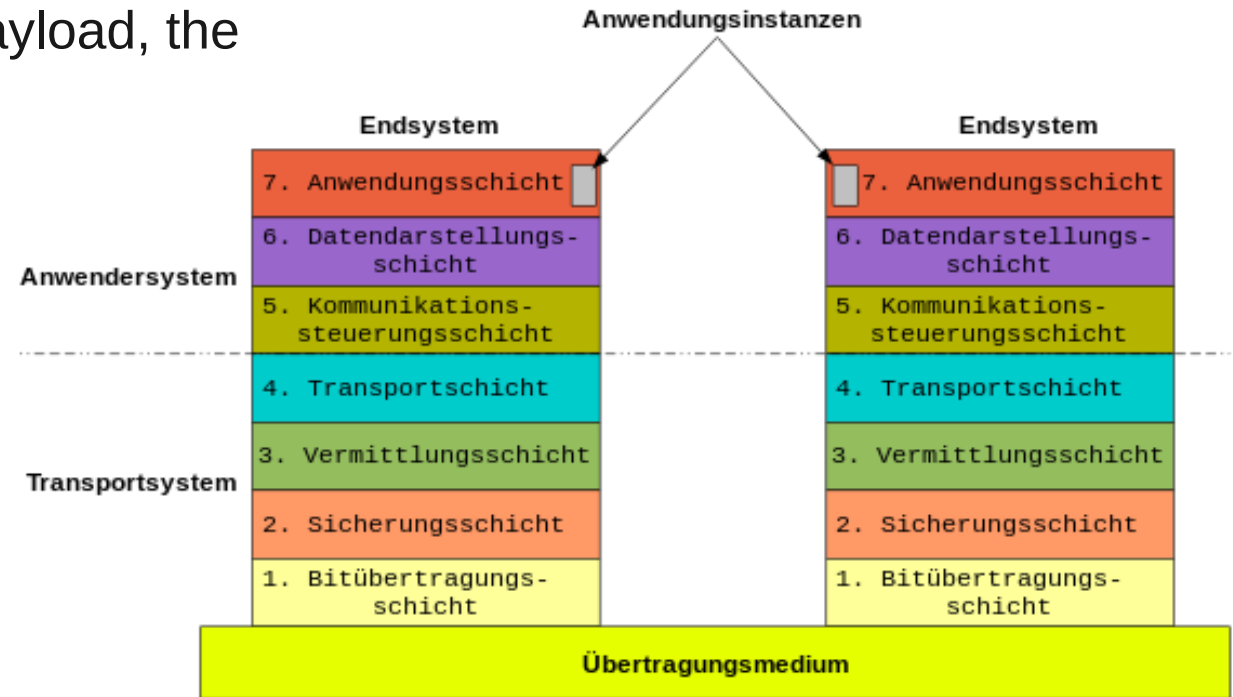
RFC 1122, Internet STD 3 (1989)								OSI model
Four layers								Seven layers
"Internet model"								OSI model
Application								Application
								Presentation
								Session
Transport								Transport
Internet								Network
Link								Data link
								Physical

Layer Abstraction

- Protocols enable an entity/instance to interact with an entity/instance at the same layer in another host
- Service definitions: provide functionality to an (N)-layer by an (N-1) layer
- Each **PDU** contains a protocol header and payload, the service data unit (**SDU**). E.g. PDU of L3:



source: https://en.wikipedia.org/wiki/OSI_model

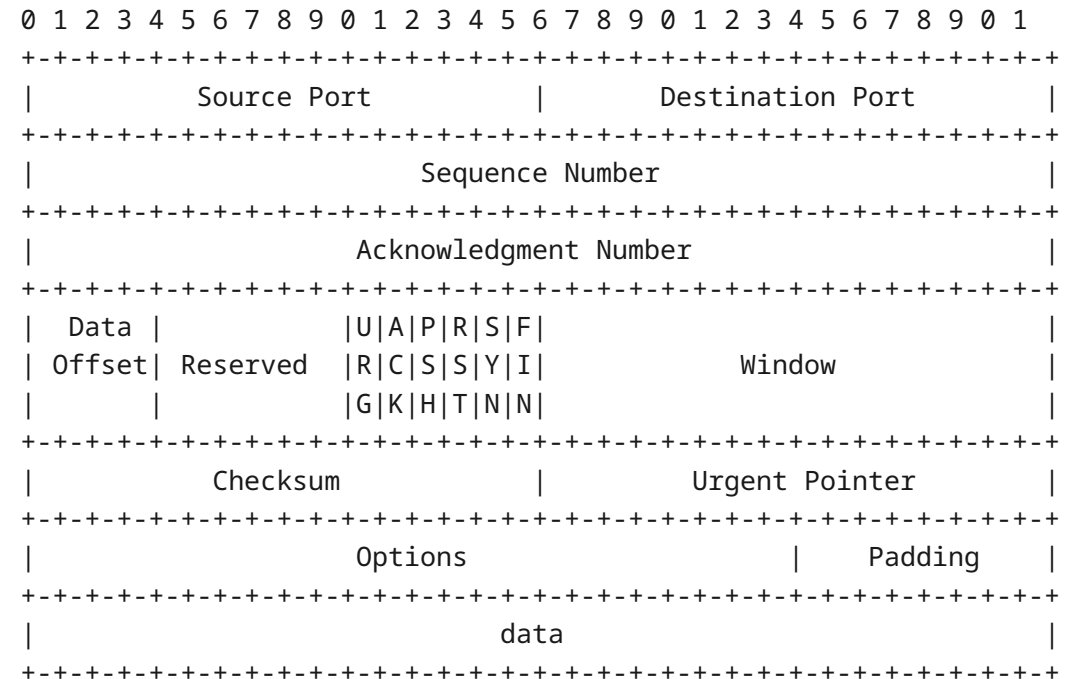


source: <https://de.wikipedia.org/wiki/OSI-Modell>



Layer 4 - Transport

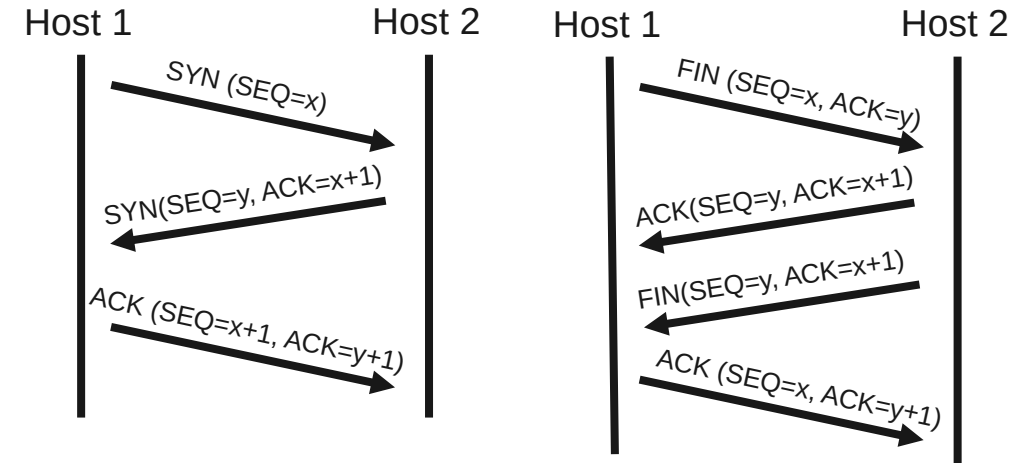
- TCP (Transmission Control Protocol)
 - Reliable (retransmission)
 - Ordered
 - Window – capacity of receiver
 - Checksum – 16bit (crc16)
 - TCP overhead: 20bytes
 - IP overhead: 20bytes
 - Ethernet frame: 18bytes (crc32)
- TCP tries to correct errors; you don't need to worry...
 - Sometimes, you need to worry...



source: <http://freesoft.org/CIE/Course/Section4/8.htm>

Layer 4 - TCP

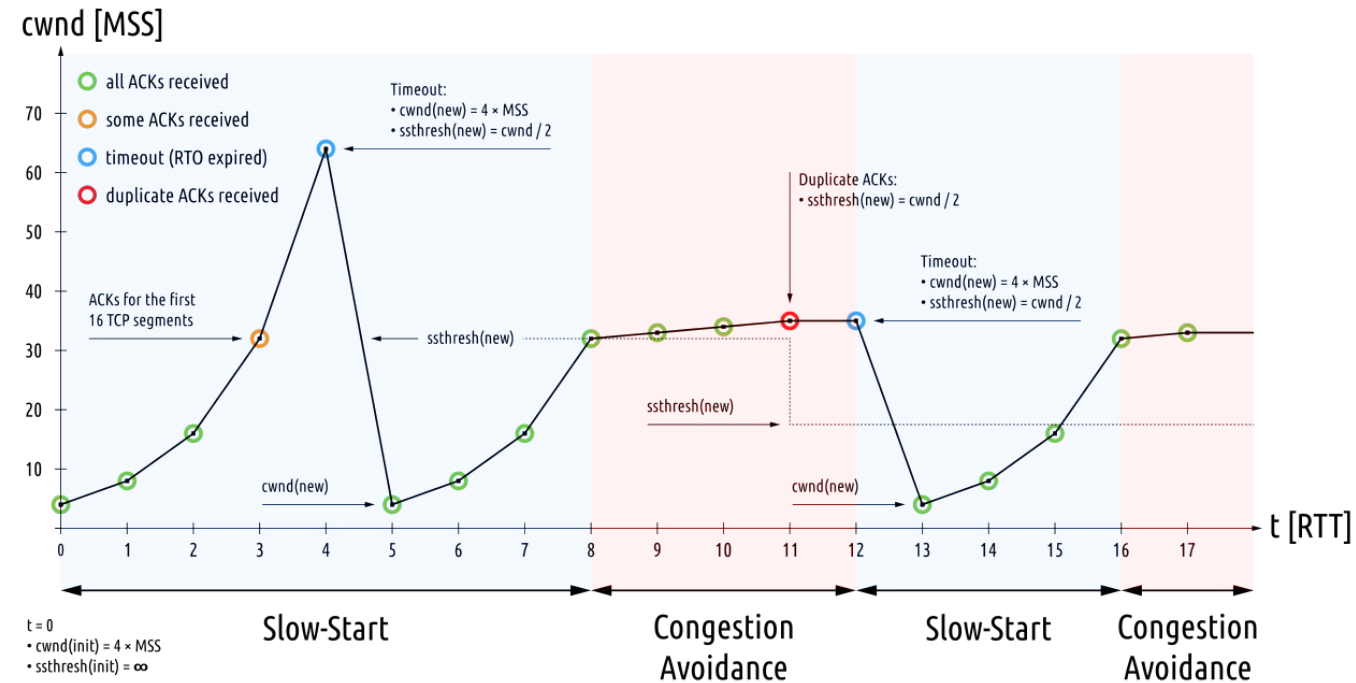
- Connection establishment
 - SYN, SYN-ACK, ACK (three way)
 - Initiates TCP session: initial sequence number is ~ random
- Connection termination
 - FIN, ACK + FIN, ACK (three/four way)
 - 3-way handshake, when host 1 sends a FIN and host 2 replies with a FIN & ACK
- Sequences and ACKs
 - Identification each byte of data
 - Order of the bytes → reconstruction
 - Detecting lost data: RTO, DupACK:



- Retransmission timeout
 - If no ACK is received after timeout (e.g. $2 \times \text{RTT}$), resend.
- Duplicate cumulative acknowledgements, selective ACK [[link](#)]
 - ACKs for last consecutive packets
 - 3 times same ACK → retransmit missing packets (fast retransmit)

Layer 4 - TCP

- Flow control
 - Sender is not overwhelming a receiver
 - Back pressure
 - Sliding window:
 - Receiver specifies the amount of additionally received data in bytes that can be buffered
 - Sender up to that amount of data before ACK
- Congestion control
 - slow-start
 - congestion avoidance
- Difference flow/congestion control



source:

https://upload.wikimedia.org/wikipedia/commons/thumb/2/24/TCP_Slow-Start_and_Congestion_Avoidance.svg/1280px-TCP_Slow-Start_and_Congestion_Avoidance.svg.png

TCP/IP from an Application Developer View

- Server in golang ([repo](#))
 - git clone
<https://github.com/tbocek/DSy>
 - Download [GoLand](#), or [others](#)
 - go run server.go → server
- Listening on TCP port 8081
 - Return string in uppercase
- Node.js version
 - Download [WebStorm](#), or [other](#)
- Client:
 - nc localhost 8081

```
const net = require('net');
const server = new net.Server();
server.listen(8081, function() {
  console.log('Launching server...');
});

server.on('connection', function(socket) {
  socket.on('data', function(chunk) {
    console.log('Data received from client: $
{chunk.toString()}');

    socket.write(chunk.toString().toUpperCase() +
"\n");
  });
});
```

```
package main
import ("bufio"
        "fmt"
        "net"
        "strings")
func main() {
  fmt.Println("Launching server...")
  ln, _ := net.Listen("tcp", ":8081") // listen
on all interfaces
  for {
    conn, _ := ln.Accept() // accept
connection on port
    message, _ :=
bufio.NewReader(conn).ReadString('\n') //read line
    fmt.Print("Message Received:",
string(message))
    newMessage := strings.ToUpper(message)
//change to upper
    conn.Write([]byte(newMessage + "\n"))
//send upper string back
  }
}
```

TCP Considerations

- Fallacy 2: Latency is zero
 - Nürnberg data center: 15ms, Australia: 300ms
 - Ping ftp.au.debian.org , [Starklink](#) + 50ms
- Problem: TCP handshake is not flexible
 - You need a handshake (1RT)
 - 1) If you want to make sure the other side accepts packets (and not drop it) - ensure both sides are ready to transmit and receive data
 - 2) If you want to exchange public / private keys
 - TCP supports 1) but not 2)
 - Use another security layer for 2), but a security layer needs at least 1 RT
 - TCP + Security = at least 2 RT
 - Nürnberg + Starlink: $2 \times (15 + 50\text{ms}) = 130\text{ms}$
 - Australian: $(2 \times 300\text{ms}) = 600\text{ms}$

- TCP + Security at least 2 RT
 - DNS query may be required too: 3 RT
 - Old security protocols add RT: 4RT
- Starklink, Kuiper, OneWeb in lower orbit
 - Viasat is 631 ms (higher orbit)
- Worst case: Starlink/Australia/DNS/TCP/old sec: 1.4s before data can be sent → new protocols on the way (HTTP/3)
- Wireshark

No.	Time	Source	Destination	Protocol	Length	Info
19	9.405192238	192.168.1.221	152.96.86.25	TLSv1.3	236	Application Data
20	9.405609773	192.168.1.221	152.96.86.25	TLSv1.3	382	Application Data
21	9.419381967	152.96.86.25	192.168.1.221	TCP	66	443 → 38598 [ACK] Seq=4565 Ack=752 Win=64640 Len=0 TSval=3439833320 TSecr=2594247482
22	9.419484927	152.96.86.25	192.168.1.221	TLSv1.3	437	Application Data
23	9.419860624	152.96.86.25	192.168.1.221	TLSv1.3	408	Application Data, Application Data
24	9.420608456	192.168.1.221	152.96.86.25	TCP	66	38598 → 443 [ACK] Seq=1068 Ack=5178 Win=64128 Len=0 TSval=2594247498 TSecr=3439833320
25	9.420685146	192.168.1.221	152.96.86.25	TLSv1.3	97	Application Data
26	9.435385870	152.96.86.25	192.168.1.221	TCP	1434	443 → 38598 [ACK] Seq=5178 Ack=1099 Win=64384 Len=1368 TSval=3439833336 TSecr=259424748
27	9.435739635	152.96.86.25	192.168.1.221	TCP	1434	443 → 38598 [ACK] Seq=6546 Ack=1099 Win=64384 Len=1368 TSval=3439833336 TSecr=259424748
28	9.435764780	192.168.1.221	152.96.86.25	TCP	66	38598 → 443 [ACK] Seq=1099 Ack=7914 Win=64128 Len=0 TSval=2594247513 TSecr=3439833336
29	9.435989541	152.96.86.25	192.168.1.221	TLSv1.3	452	Application Data, Application Data
30	9.476608117	192.168.1.221	152.96.86.25	TCP	66	38598 → 443 [ACK] Seq=1099 Ack=8300 Win=64128 Len=0 TSval=2594247554 TSecr=3439833336
1419	68.024803369	192.168.1.221	152.96.86.25	TLSv1.3	105	Application Data
1420	68.038496928	152.96.86.25	192.168.1.221	TLSv1.3	105	Application Data
1421	68.038539673	192.168.1.221	152.96.86.25	TCP	66	38598 → 443 [ACK] Seq=1138 Ack=8339 Win=64128 Len=0 TSval=2594306116 TSecr=3439891939



Wireshark – sometimes needed when designing protocols

- Decrypt TLS
- export SSLKEYLOGFILE=/tmp/keylogfile.txt

The screenshot displays the Wireshark interface for a network capture titled '*enx106530e2c54d'. The filter bar shows the filter 'ip.src==152.96.80.48 || ip.dst==152.96.80.48'. The packet list pane shows a series of packets, with packet 30 selected. The packet details pane shows the following information for the selected packet:

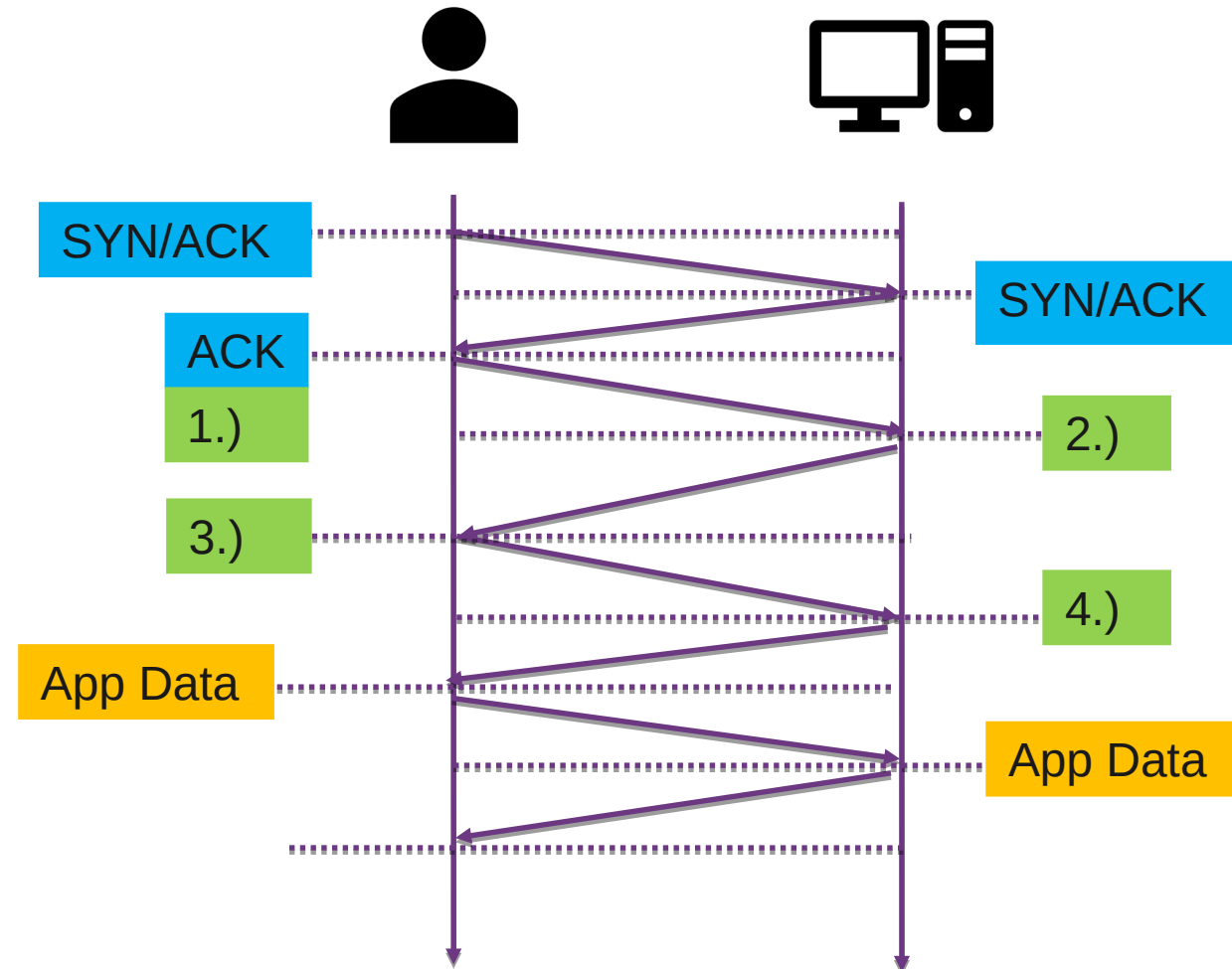
```
> Transmission Control Protocol, Src Port: 50908, Dst Port: 443, Seq: 0, Len: 0
  Source Port: 50908
  Destination Port: 443
  [Stream index: 5]
  [TCP Segment Len: 0]
  Sequence number: 0 (relative sequence number)
  [Next sequence number: 0 (relative sequence number)]
  Acknowledgment number: 0
  1010 ... = Header Length: 40 bytes (10)
  - Flags: 0x0c2 (SYN, ECN, CWR)
    000. .... = Reserved: Not set
```

The packet bytes pane shows the raw data of the packet:

```
0000 00 08 e3 ff fd 90 10 65 30 e2 c5 4d 08 00 45 00 .....e 0...M..E.
0010 00 3c 2f d6 40 00 40 06 b3 a0 98 60 d6 54 98 60 </@.@...T..
0020 50 30 c6 dc 01 bb 7d 4a 53 f3 00 00 00 00 a0 c2 P0...}J S.....
0030 72 10 57 74 00 00 02 04 05 b4 04 02 08 0a ba 62 r.Wt.....b
0040 7d 59 00 00 00 00 01 03 03 07 }Y.....
```

Layer 4 – TCP + TLS

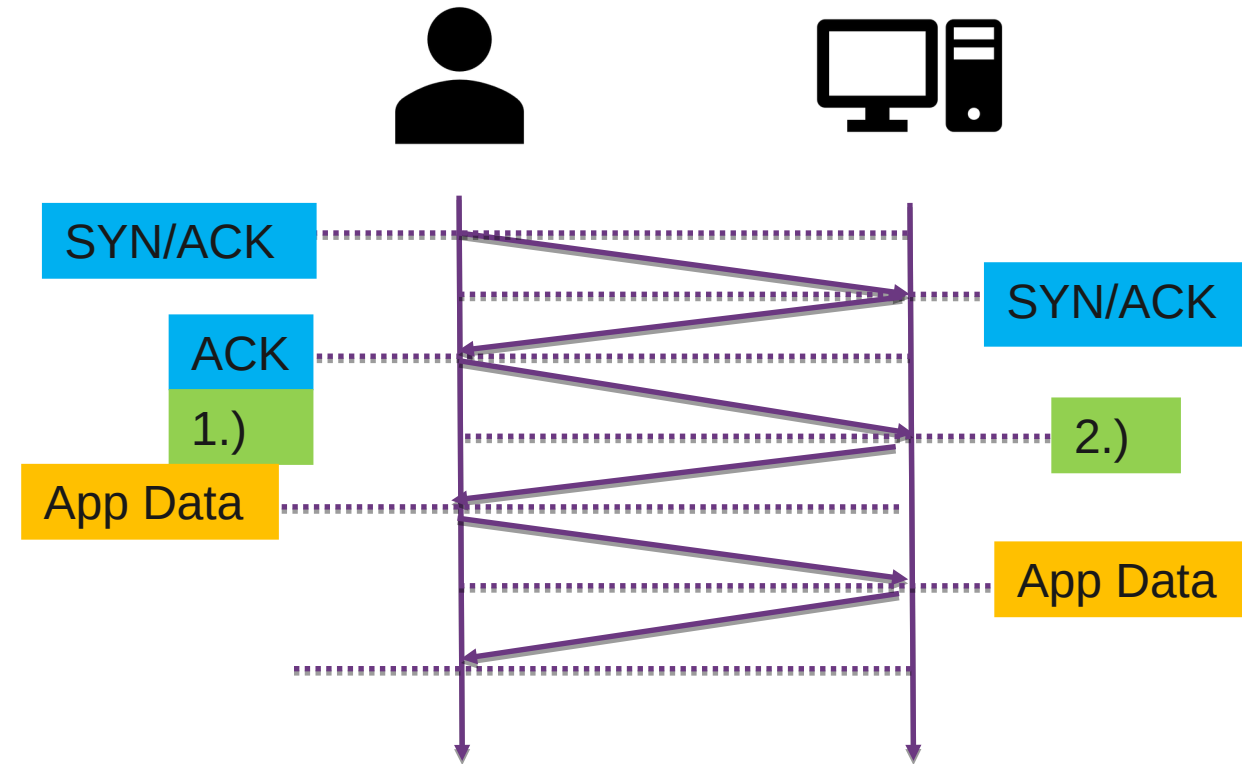
- Security: Transport Layer Security (TLS)
 1. "client hello" lists cryptographic information, TLS version, ciphers/keys
 2. "server hello" chosen cipher, the session ID, random bytes, digital certificate (checked by client), optional: "client certificate request"
 3. Key exchange using random bytes, now server and client can calc secret key
 4. "finished" message, encrypted with the secret key
- 3 RTT to send first byte, 4RTT to receive first byte



```
PING sydney.edu.au (129.78.5.8) 56(84) bytes of data.  
64 bytes from scilearn.sydney.edu.au (129.78.5.8): icmp_seq=1 ttl=233 time=307 ms  
64 bytes from scilearn.sydney.edu.au (129.78.5.8): icmp_seq=2 ttl=233 time=305 ms  
64 bytes from scilearn.sydney.edu.au (129.78.5.8): icmp_seq=3 ttl=233 time=305 ms
```

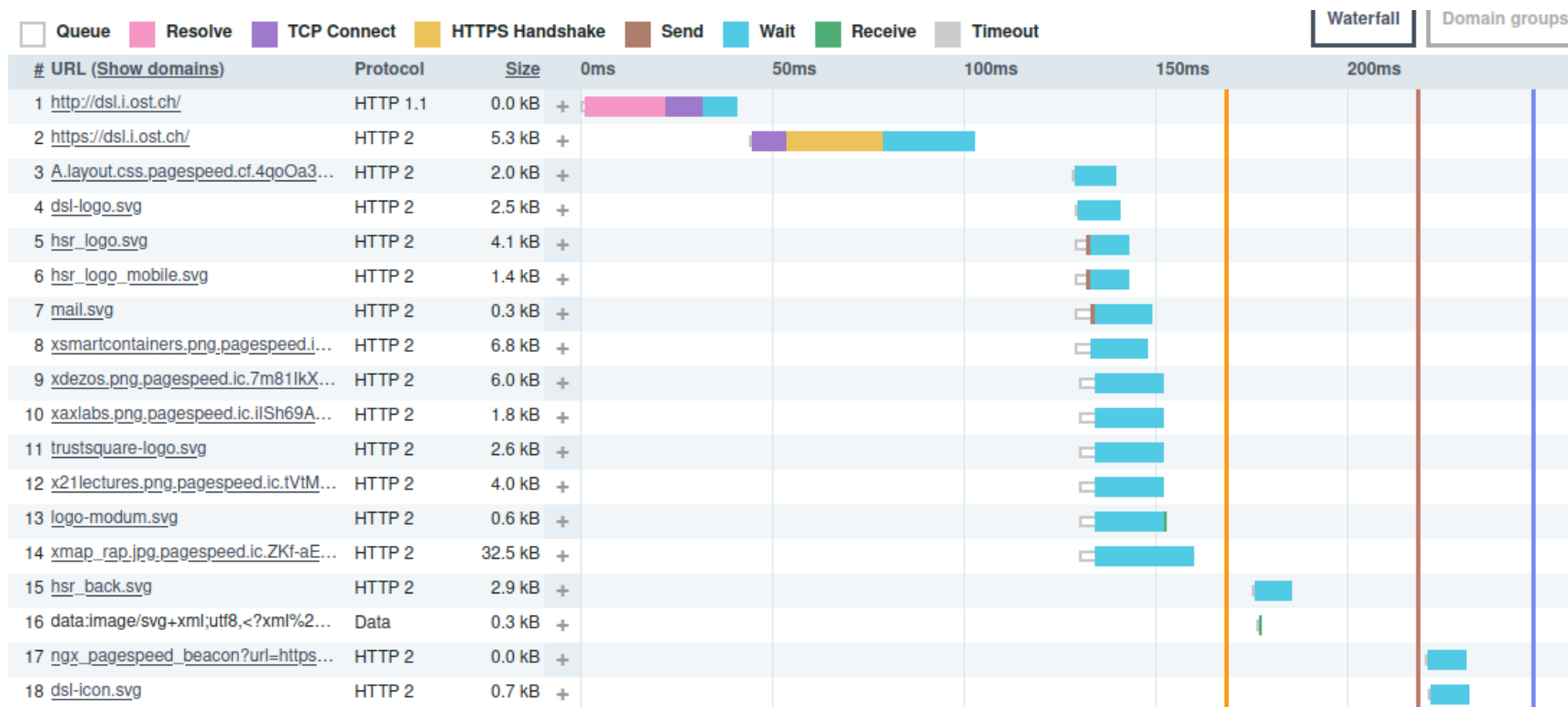
Layer 4 – TCP + TLS

- TCP + TLS handshake:
 - 1RTT - Ping to Australia: 329ms
 - 3RTT = 987ms! No data sent yet
- TLS 1.3, finished Aug 2018
 - 1 RTT instead of 2
 - 1.) Client Hello, Key Share
 - 2.) Server Hello, key Share, Verify Certificate, Finished
 - 0 RTT possible, for previous connections, loosing perfect forward secrecy
- 95% of browsers used already support it



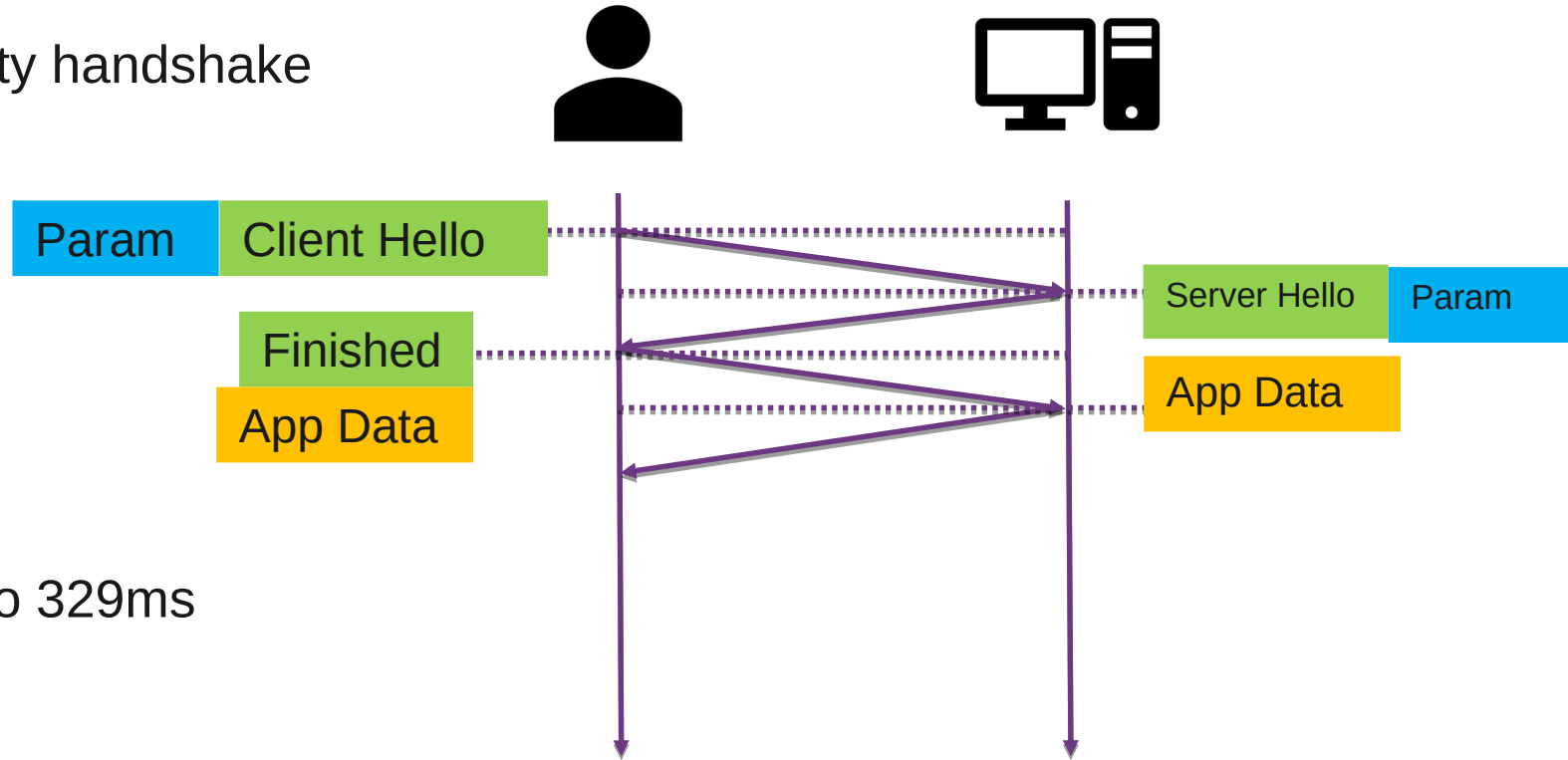
Layer 4 – TCP + TLS

- Website Speed Test [[link](#)]
 - Resolve → DNS, TCP Connect → TCP Handshake, HTTPS Handshake → TLS/SSL



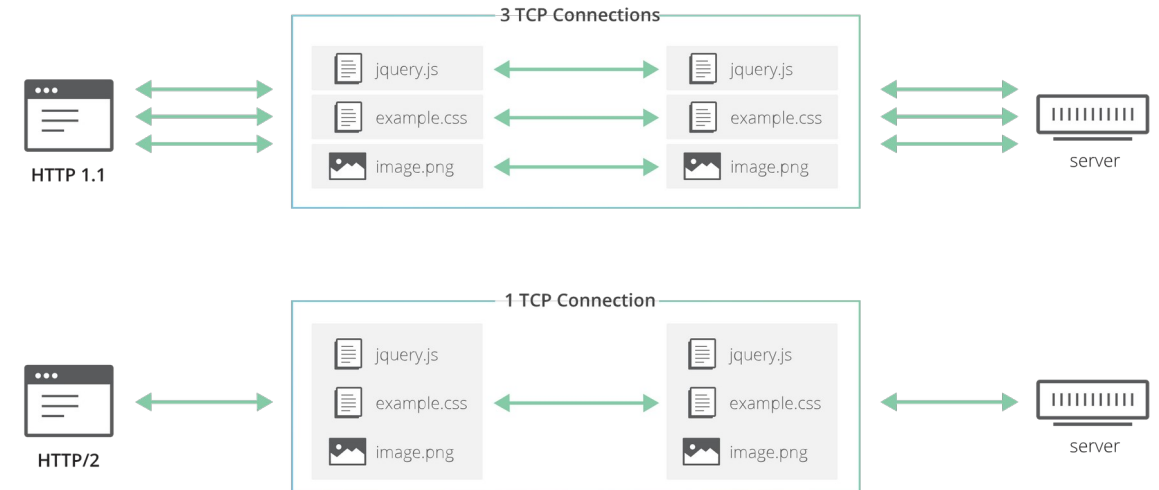
QUIC / HTTP3

- QUIC: 1RTT connection + security handshake
 - For known connections: 0RTT
 - Built in security
 - Reports
 - Facebook
 - state of HTTP
- Example Australia: from 987ms to 329ms

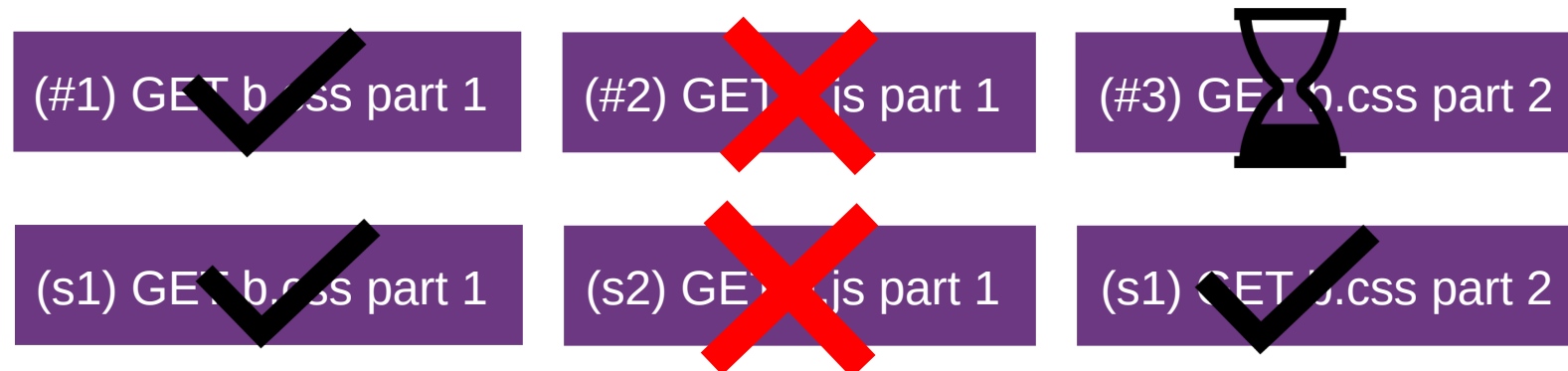


QUIC / HTTP3

- Multiplexing in HTTP/2
 - [HTTP/1 → HTTP/2](#)
- HTTP/2: Head-of-line blocking
 - One packet loss, TCP needs to be ordered
 - QUIC can multiplex requests: one stream does not affect others
- HTTP/3 is great, but...
 - NAT → SYN, ACK, FIN, conntrack knows when connection ends, not with QUIC, timeouts, new entries, many entries
 - HTTP header compression, referencing previous headers
 - Many TCP [optimizations](#)

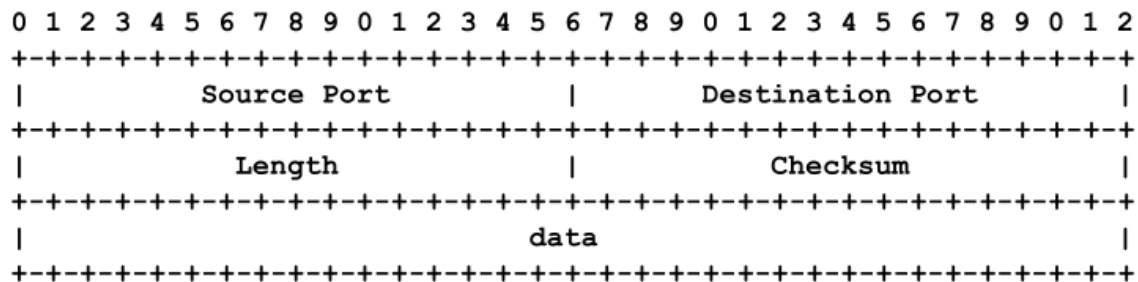


source: <https://blog.cloudflare.com/the-road-to-quic/>



Layer 4 - Transport

- User Datagram Protocol (UDP)
 - UDP is used for DNS, streaming audio and video
 - Simple connectionless communication model
 - No guarantee
 - Delivery
 - Ordering
 - Duplicate protection



- SCTP (Stream Control Transmission Protocol)
 - Message-based
 - Allows data to be divided into multiple streams
 - Syn cookies - SCTP uses a four-way handshake with a signed cookie.
 - Multi-homing multiple IP addresses of endpoints
 - Not widely used: “We have been deploying SCTP in several applications now, and encountered significant problem with SCTP support in various home routers.”
 - E.g., OpenWRT – not enabled by default
 - E.g., UFW - Uncomplicated Firewall – not supported
 - SCTP used by WebRTC, but tunneled over UDP

UDP example

- UDP Server (Java)

```
import java.net.*;

class Server
{
    public static void main(String args[]) throws
Exception
    {
        DatagramSocket serverSocket = new
DatagramSocket(8081);
        byte[] receiveData = new byte[1024];
        while(true)
        {
            DatagramPacket receivePacket = new
DatagramPacket(receiveData, receiveData.length);
            serverSocket.receive(receivePacket);
            String s = new
String( receivePacket.getData());
            System.out.println("Message Received: " +
s);
        }
    }
}
```

- UDP Client (golang)

```
package main

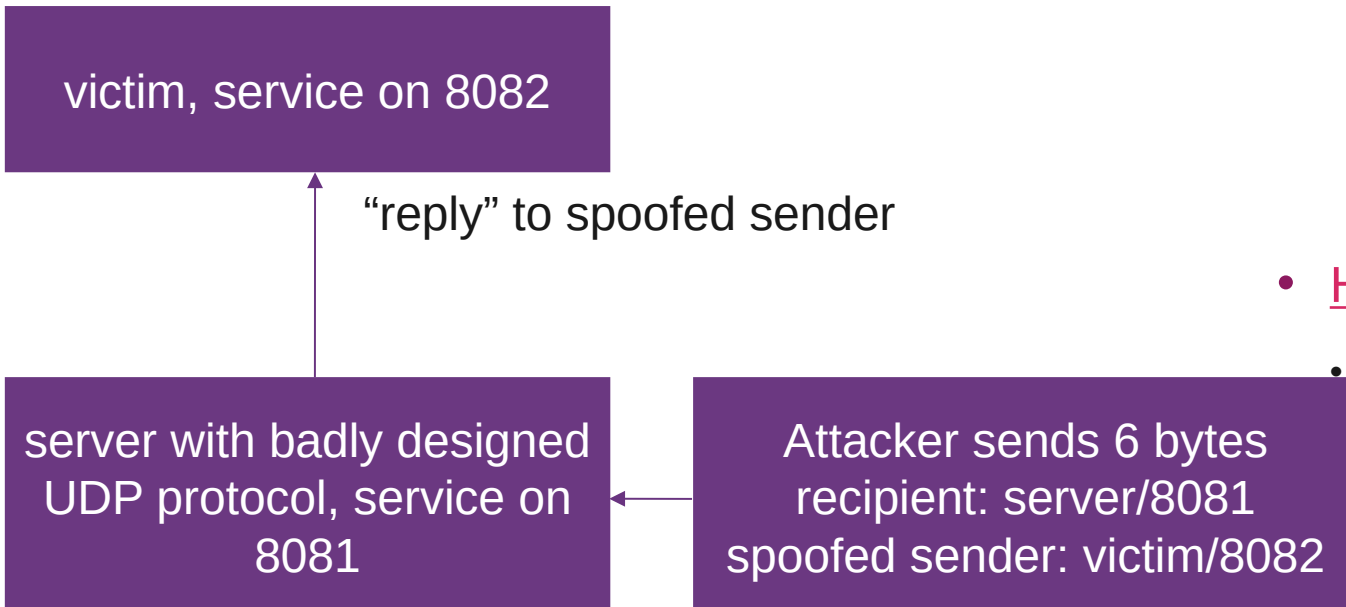
import (
    "net"
)

func main() {
    srv, _ :=
net.ResolveUDPAddr("udp", "127.0.0.1:8081")
    local, _ := net.ResolveUDPAddr("udp",
"127.0.0.1:0")
    conn, _ := net.DialUDP("udp", local, srv)
    defer conn.Close()
    conn.Write([]byte("5Anybody there?"))
}
```

- nc -u localhost 8081

Layer 4 - Transport

- DDoS Amplification Attacks
 - Request 10 bytes, reply 100 bytes → factor 10
- Local demo with server-ra/victim, and hping3
 - `hping3 --udp IP -p 8081 -E test.tmp -d 6 -s 8082 -c 1`



- Attacker in go/Java/node/c#
 - You need to spoof UDP packets, typically not supported in those languages
 - Go: `func DialUDP(network string, laddr, raddr *UDPAddr) (*UDPConn, error)`
 - laddr: we need to set here the victims IP/port
 - But go tries to bind to that
 - Not yours: “bind: cannot assign requested address”
- Hping3: Pen test tool
 - hping3 is a command-line oriented IP, TCP, UDP, ICMP and RAW-IP packet assembler

Comparison – Transport Layer

TCP *

UDP *

SCTP *

(QUIC) *

- | | | | |
|--------------------------------|-------------------------------|-------------------------------|-----------------------|
| ■ Transport layer | ■ Transport layer | ■ Transport layer | ■ Transport layer* |
| ■ Connection oriented | ■ Connection less | ■ Connection oriented | ■ Connection oriented |
| ■ Reliable transfer | ■ Unreliable transfer | ■ Reliable transfer | ■ Reliable transfer |
| ■ Streams | ■ Messages | ■ Messages | ■ Multistream |
| ■ Guaranteed order | ■ Unordered | ■ User can choose | ■ Guaranteed order |
| ■ Widely used – HTTP/1, HTTP/2 | ■ Widely used – DNS, HTTP/3 | ■ WebRTC | ■ HTTP/3 |
| ■ Flow and congestion control | ■ No flow, congestion | ■ Flow and congestion control | ■ Flow and congestion |
| ■ Heavyweight | ■ Lightweight | ■ Heavyweight | ■ Heavyweight* |
| ■ Error checking and recovery | ■ Error checking, no recovery | ■ Error checking and recovery | ■ Integrity check |