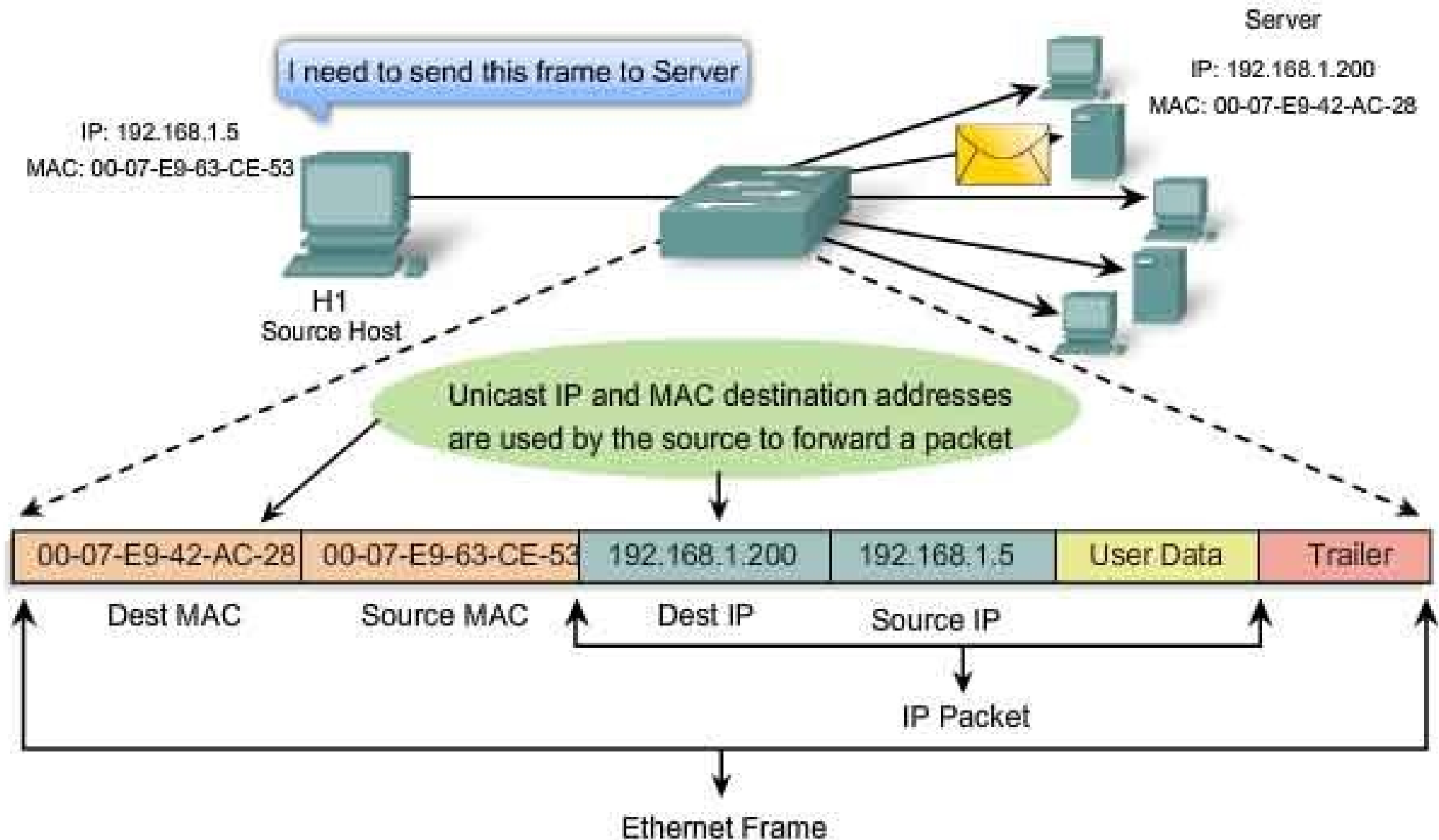# 7 – Protocols

**Marian Marinov**
**CEO of 1H Ltd.**
**mm@1h.com**

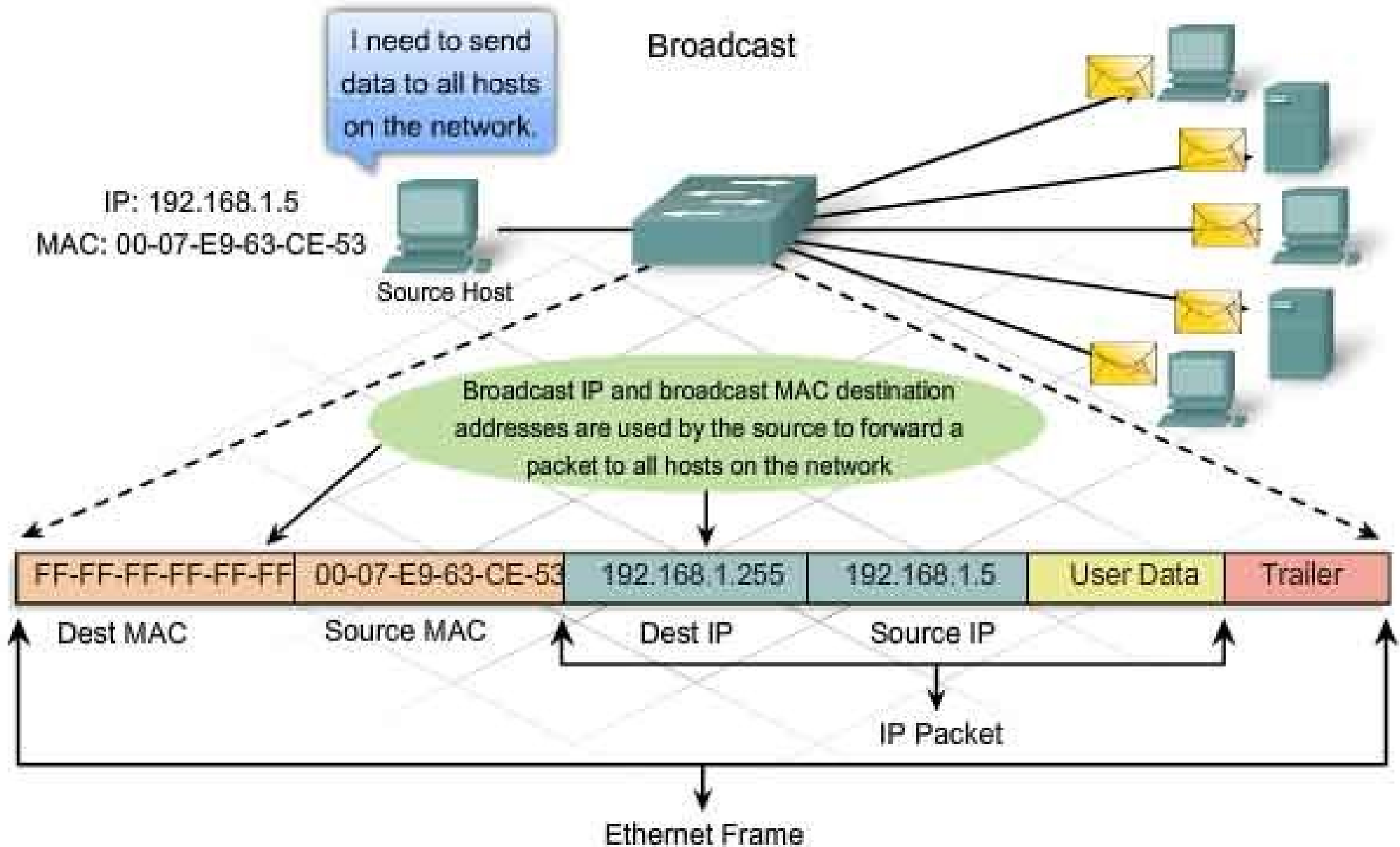**Borislav Varadinov**
**System Administrator**
**bobi [ at ] itp.bg**

- **ARP/RARP**
- **ICMP**
- **UDP**
- **TCP**
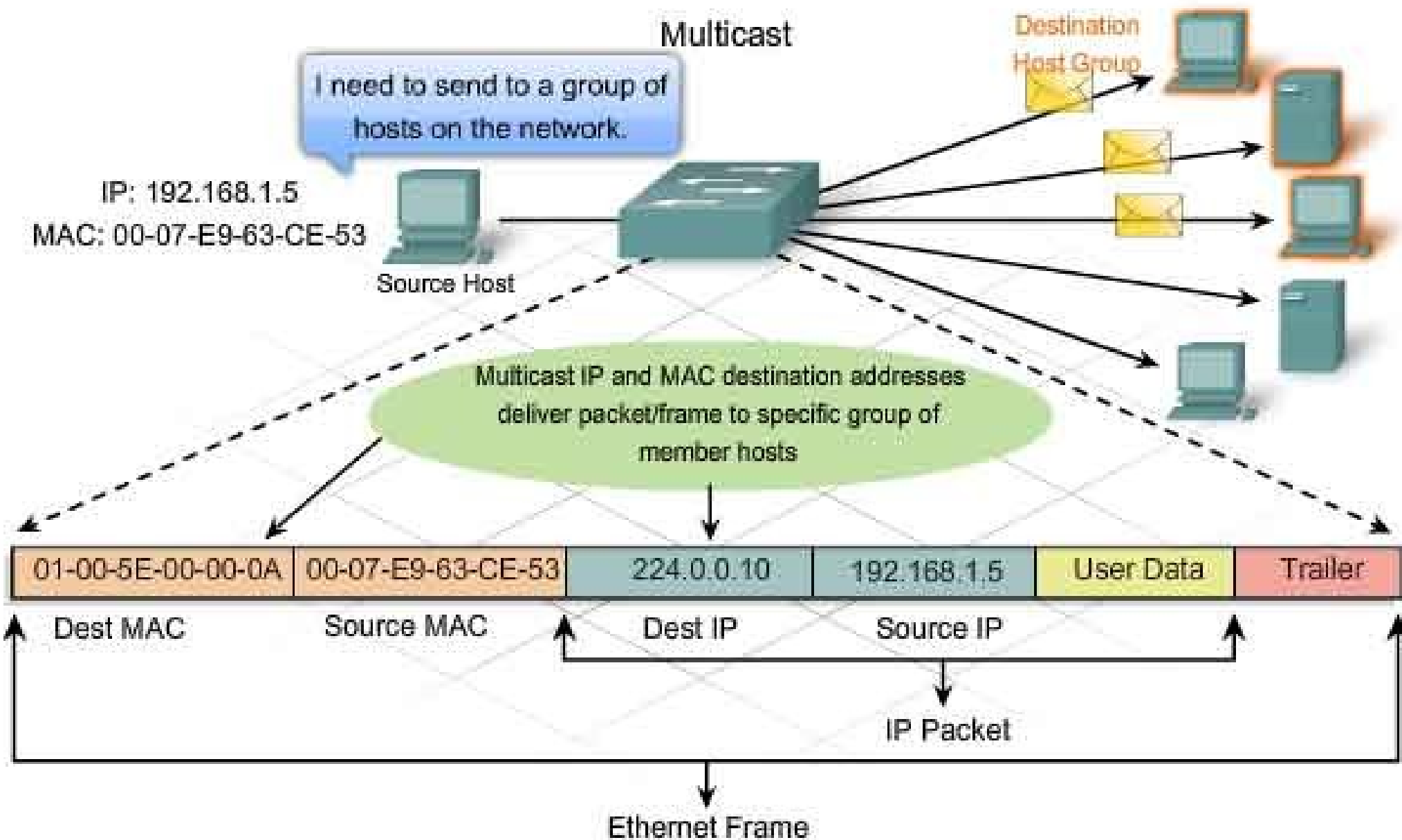- **TCP Congestion**
- **SCTP**
- **DCCP**
- **DNS**

# Type of requests - Unicast

# Type of requests - Broadcast

# Type of requests - Multicast

# Address Resolution Protocol

- **Address resolution**
  - **Forward**
  - **Reverse**

- **ARP**
  - Probe
  - Proxy
  - Mediation
  - Stuffing

# Address Resolution Protocol

- **Address resolution**

    - **Forward** (what is the MAC of this machine)

**Request**

ARP header

| 0 | 7 | 15 | 31 |
|---|---|---|---|

| Hardware type | | Protocol type **0x0800** | |
|---|---|---|---|
| Hardware address length | Protocol address length | **1 - req** Opcode **2 - reply** | |
| **08:11:96:03:B2:28** Source hardware address | | | |
| **192.168.2.254** Source protocol address | | | |
| **FF:FF:FF:FF:FF:FF** Destination hardware address | | | |
| **192.168.2.58** Destination protocol address | | | |

# Address Resolution Protocol

- **Address resolution**

    – **Forward** (what is the MAC of this machine)

**Reply**

ARP header

| 0 | 7 | 15 | 31 |
|---|---|---|---|

| Hardware type | | Protocol type | |
|---|---|---|---|
| Hardware address length | Protocol address length | Opcode | |
| **40:b3:95:80:c5:aa** | | Source hardware address | |
| **192.168.2.58** | | Source protocol address | |
| **08:11:96:03:b2:28** | | Destination hardware address | |
| **192.168.2.254** | | Destination protocol address | |

# Address Resolution Protocol

- **Address resolution**
  - **Reverse** (what is the IP of this machine)

**Request**

ARP header

| 0 | 7 | 15 | 31 |
|---|---|---|---|
| Hardware type | | Protocol type | |
| Hardware address length | Protocol address length | Opcode | |
| 08:11:96:03:B2:28 | Source hardware address | | |
| 192.168.2.254 | Source protocol address | | |
| 40:b3:95:80:c5:aa | Destination hardware address | | |
| 0.0.0.0 | Destination protocol address | | |

# Address Resolution Protocol

- **Address resolution**

  - **Reverse** (what is the IP of this machine)

**Reply**

ARP header

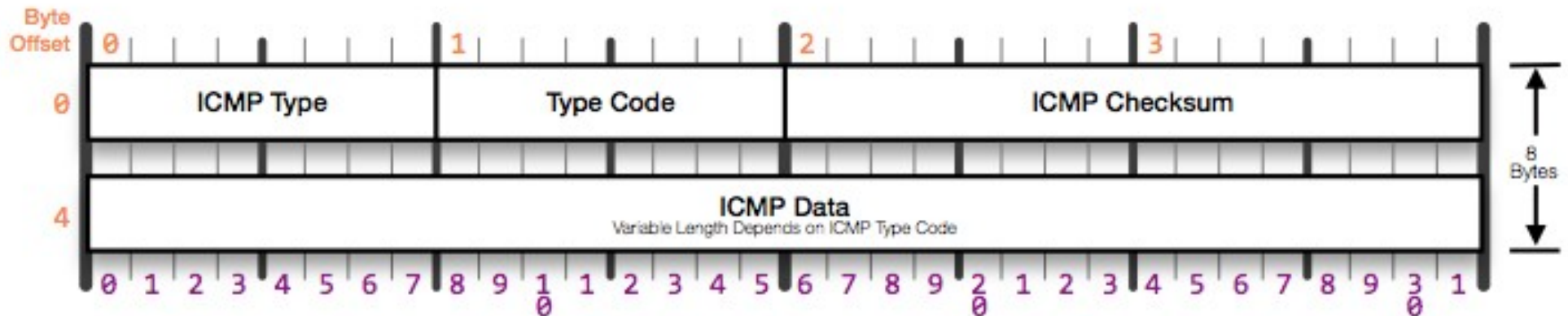| 0 | 7 | 15 | 31 |
|---|---|---|---|
| Hardware type | | Protocol type | |
| Hardware address length | Protocol address length | Opcode | |
| **40:b3:95:80:c5:aa** | | Source hardware address | |
| **192.168.2.58** | | Source protocol address | |
| **08:11:96:03:B2:28** | | Destination hardware address | |
| **192.168.2.254** | | Destination protocol address | |

# Address Resolution Protocol

- **How it actually looks**

15:12:43.772954 ARP, Ethernet (len 6), IPv4 (len 4),
Request who-has 192.168.2.58 tell 192.168.2.254, length 28
15:12:43.962834 ARP, Ethernet (len 6), IPv4 (len 4),
Reply 192.168.2.58 is-at 40:b3:95:80:c5:aa, length 46

- ARP probe
- ARP proxy
- ARP mediation
- ARP stuffing

# ICMP Header
RFC 792 Outlines the ICMP Protocol

Byte Offset

| 0 | | | 1 | | | 2 | | | 3 | | |

| 0 | ICMP Type | Type Code | ICMP Checksum |

8 Bytes

| 4 | ICMP Data Variable Length Depends on ICMP Type Code |

0 1 2 3 4 5 6 7 8 9 10 1 2 3 4 5 6 7 8 9 20 1 2 3 4 5 6 7 8 9 30 1

**ICMP Type**
0  Echo Reply

**ICMP Type**
3  Destination Unreachable

Type Code
- 0  Network Unreachable
- 1  Host Unreachable
- 2  Protocol Unreachable
- 3  Port Unreachable
- 4  Fragment Necessary
- 5  Source Route Failed
- 6  Destination Network Unknown
- 7  Destination Host Unknown
- 8  Obsolete
- 9  Destination Network Prohibited
- 10  Destination Host Prohibited
- 11  Network Unreachable for TOS
- 12  Host Unreachable for TOS
- 13  Communication Prohibited

**ICMP Type**
4  Source Quench

**ICMP Type**
5  Redirect

Type Code
- 0  Redirect for Network
- 1  Redirect for Host
- 2  Redirect for TOS and Network
- 3  Redirect for TOS and Host

**ICMP Type**
8  Echo Request

**ICMP Type**
9  Router Advertisement

**ICMP Type**
10  Router Solicitation

**ICMP Type**
11  Time to Live Exceeded

Type Code
- 0  TTL Exceeded in Transit
- 1  TTL Exceeded in Reassembly

**ICMP Type**
12  Parameter Problem

Type Code
- 0  Pointer Problem
- 1  Required Option Missing

**ICMP Type**
13  Timestamp Request

**ICMP Type**
14  Timestamp Reply

**ICMP Type**
17  Address Mask Request

**ICMP Type**
18  Address Mask Reply

**ICMP QUERY OR RESPONSE**

**ICMP ERROR MESSAGE**

ICMP Protocol Header Format
Created by Troy Jessup - http://www.troyjessup.com

Telerik Academy

# Internet Control Message Protocol - ICMP

- **ICMP types**
  - 0 – Echo replay
  - 1,2 – Reserved
  - 3 – Destination unreachable
  - 8 – Echo request
  - 9 – TTL Exceeded
  - 30 – Traceroute

# Internet Control Message Protocol - ICMP

- **Type codes**

  **11 – Time to live exceeded**
  **0 – in transit**
  **1 – in reassembly**

  **3 – Destination unreachable**
  **0 – network unreachable**
  **1 – host unreachable**
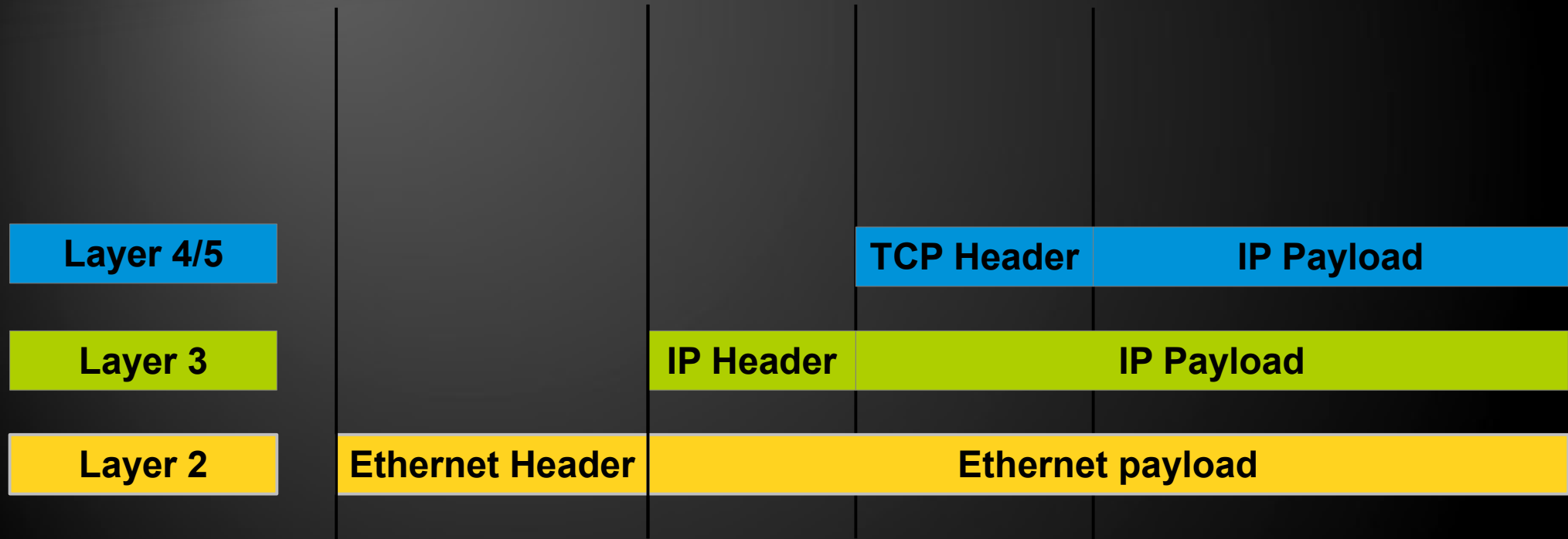  **2 – protocol unreachable**
  **3 – port unreachable**
  **6 – network unknown**
  **7 – host unknown**
  **9 – network prohibited**
  **10 – host prohibited**

# Protocol Encapsolation

| Layer 4/5 | | | | TCP Header | IP Payload |
| Layer 3 | | | IP Header | IP Payload | |
| Layer 2 | Ethernet Header | Ethernet payload | | | |

User Datagram Protocol - UDP

Machine X

Machine Y

UDP (RFC768 Jon Postel 1980)

# User Datagram Protocol - UDP

Dated 12:59 PM 07/01/2010

## UDP Header – RFC 768

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| Source Port | Destination Port |
|-------------|------------------|
| 0 · 1 | 2 · 3 |
| Length | Checksum |
| 4 · 5 | 6 · 7 |
| Data | |
| 8 · 9 | 10 · 11 |

### Common UDP Well-Known Ports

| Port | Description |
|------|-------------|
| 7 | Echo |
| 19 | Chargen |
| 37 | Time |
| 53 | Domain |
| 67 | Bootps (DHCP) |
| 68 | Bootpc (DHCP) |
| 69 | Tftp |
| 137 | Netbios-ns |

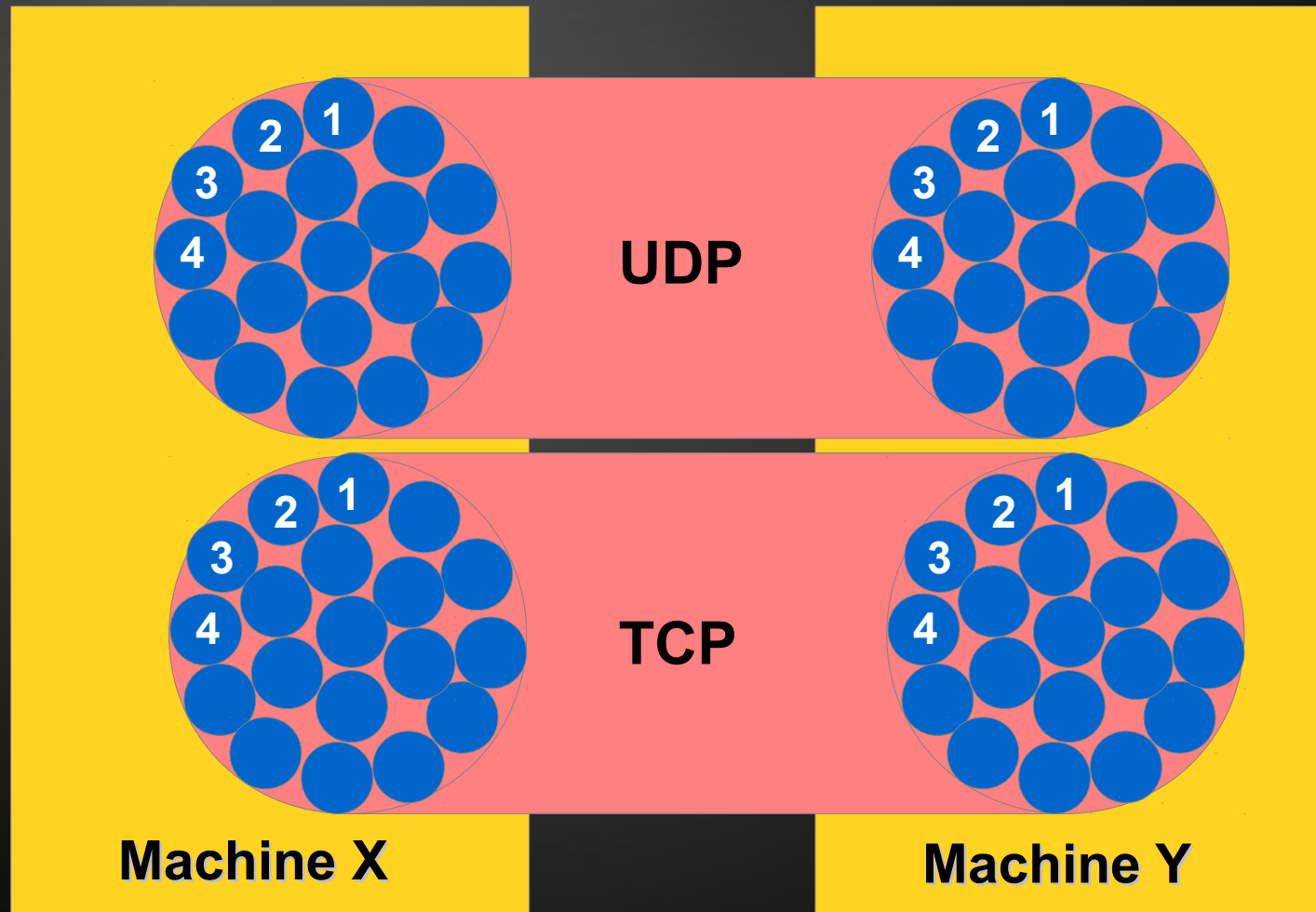| Port | Description |
|-------|-------------|
| 138 | Netbios-dgm |
| 161 | Snmp |
| 162 | Snmp-trap |
| 500 | Isakmp |
| 514 | Syslog |
| 520 | Rip |
| 33434 | Traceroute |

**Length**
The number of bytes in the entire datagram, including the header; minimum value = 8

**Checksum**
Calculated using a pseudo header that includes the IP source and destination addresses, protocol and UDP length, UDP header and data.
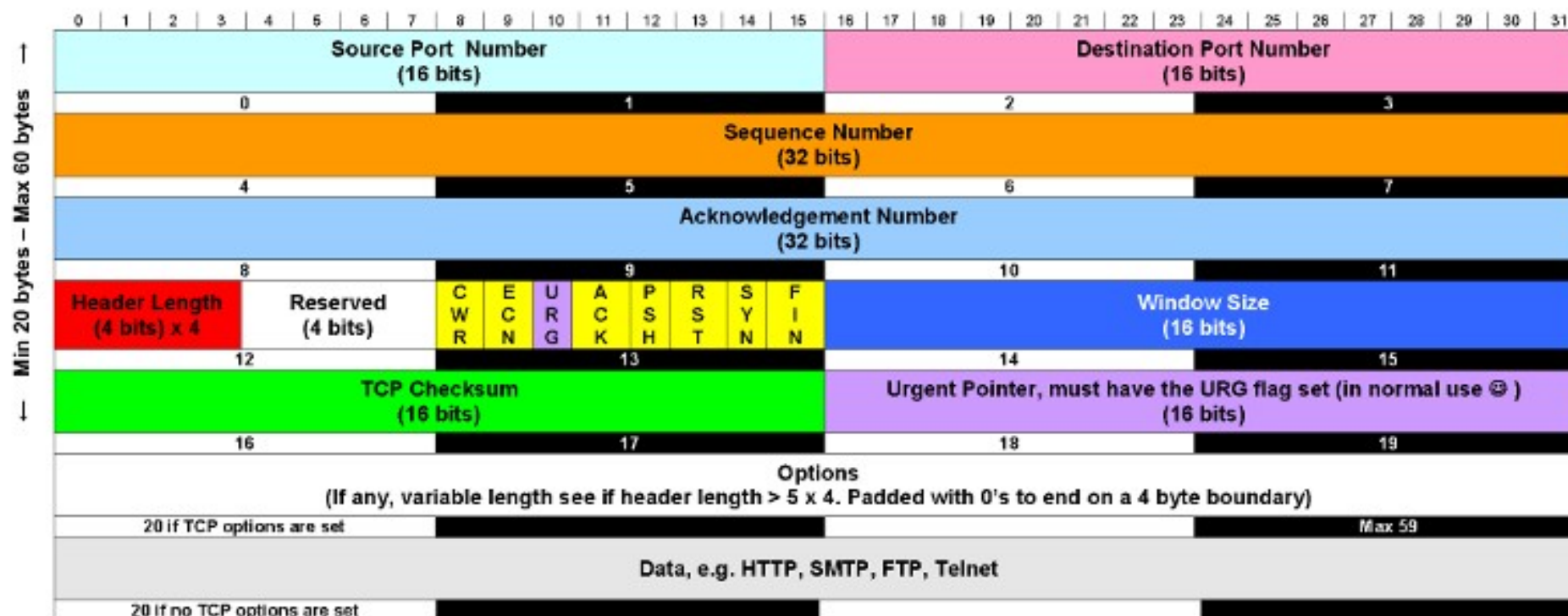
# Transmission Datagram Protocol - TCP

**UDP**

**TCP**

**Machine X**

**Machine Y**

**TCP (RFC793 Jon Postel 1981)**

# Transmission Datagram Protocol - TCP

Dated 1:00 PM 07/01/2010

## TCP Header – RFC 793

Min 20 bytes – Max 60 bytes
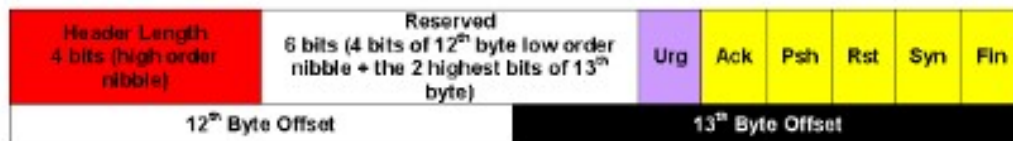
| 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 | 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 |
|---|---|
| **Source Port Number** (16 bits) | **Destination Port Number** (16 bits) |
| 0 | 1 | 2 | 3 |

| **Sequence Number** (32 bits) |
|---|
| 4 | 5 | 6 | 7 |

| **Acknowledgement Number** (32 bits) |
|---|
| 8 | 9 | 10 | 11 |

| **Header Length** (4 bits) x 4 | **Reserved** (4 bits) | C W R | E C N | U R G | A C K | P S H | R S T | S Y N | F I N | **Window Size** (16 bits) |
|---|---|---|---|---|---|---|---|---|---|---|
| 12 | 13 | 14 | 15 |

| **TCP Checksum** (16 bits) | Urgent Pointer, must have the URG flag set (in normal use ☺ ) (16 bits) |
|---|---|
| 16 | 17 | 18 | 19 |

| Options (If any, variable length see if header length > 5 x 4. Padded with 0's to end on a 4 byte boundary) |
|---|
| 20 if TCP options are set ... Max 59 |

| Data, e.g. HTTP, SMTP, FTP, Telnet |
|---|
| 20 if no TCP options are set |

Old TCP Flags (13th Byte Offset)

| Reserved | Urg | Ack | Psh | Rst | Syn | Fin |
|---|---|---|---|---|---|---|

Note on the new flags :-
CWR – Congestion Window Reduced
ECN – Explicit Congestion Notice

| Header Length 4 bits (high order nibble) | Reserved 6 bits (4 bits of 12th byte low order nibble + the 2 highest bits of 13th byte) | Urg | Ack | Psh | Rst | Syn | Fin |
|---|---|---|---|---|---|---|---|
| 12th Byte Offset | 13th Byte Offset | | | | | | |

Notes :-
If the header length has 0x05 which is 20 in the real world, the TCP options are present. Other than the initial SYN all other communications should have the ACK flag set. If the Urg flag is set the packet may contain control or interrupt characters.
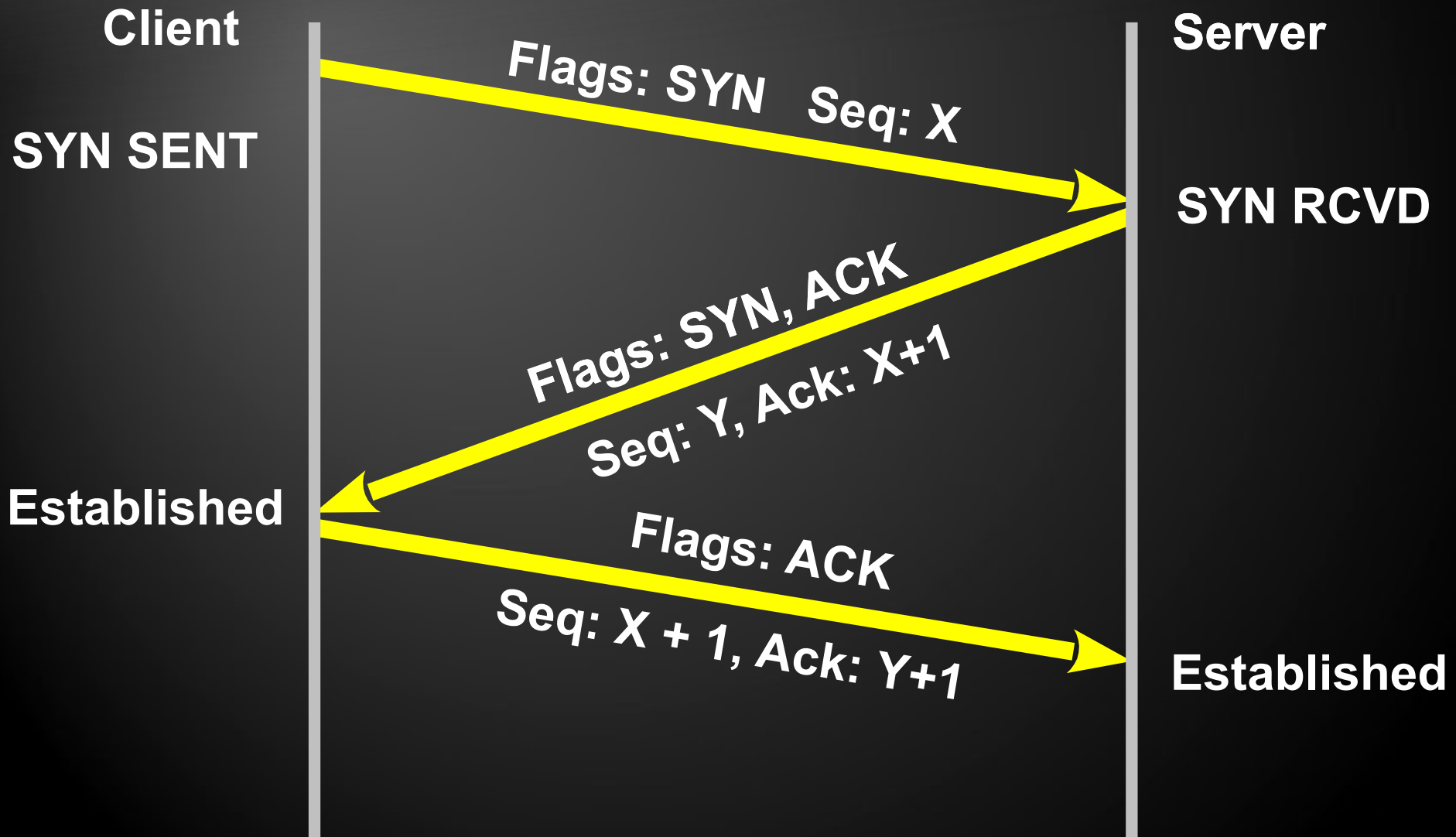
# Transmission Datagram Protocol - TCP

- ## TCP (RFC793 Jon Postel 1981)
  - Session establishment and tear-down
  - Window procedure
  - Slow start and congestion avoidance (Van Jacobson 1988)
  - Fast open
  - Syn cookies

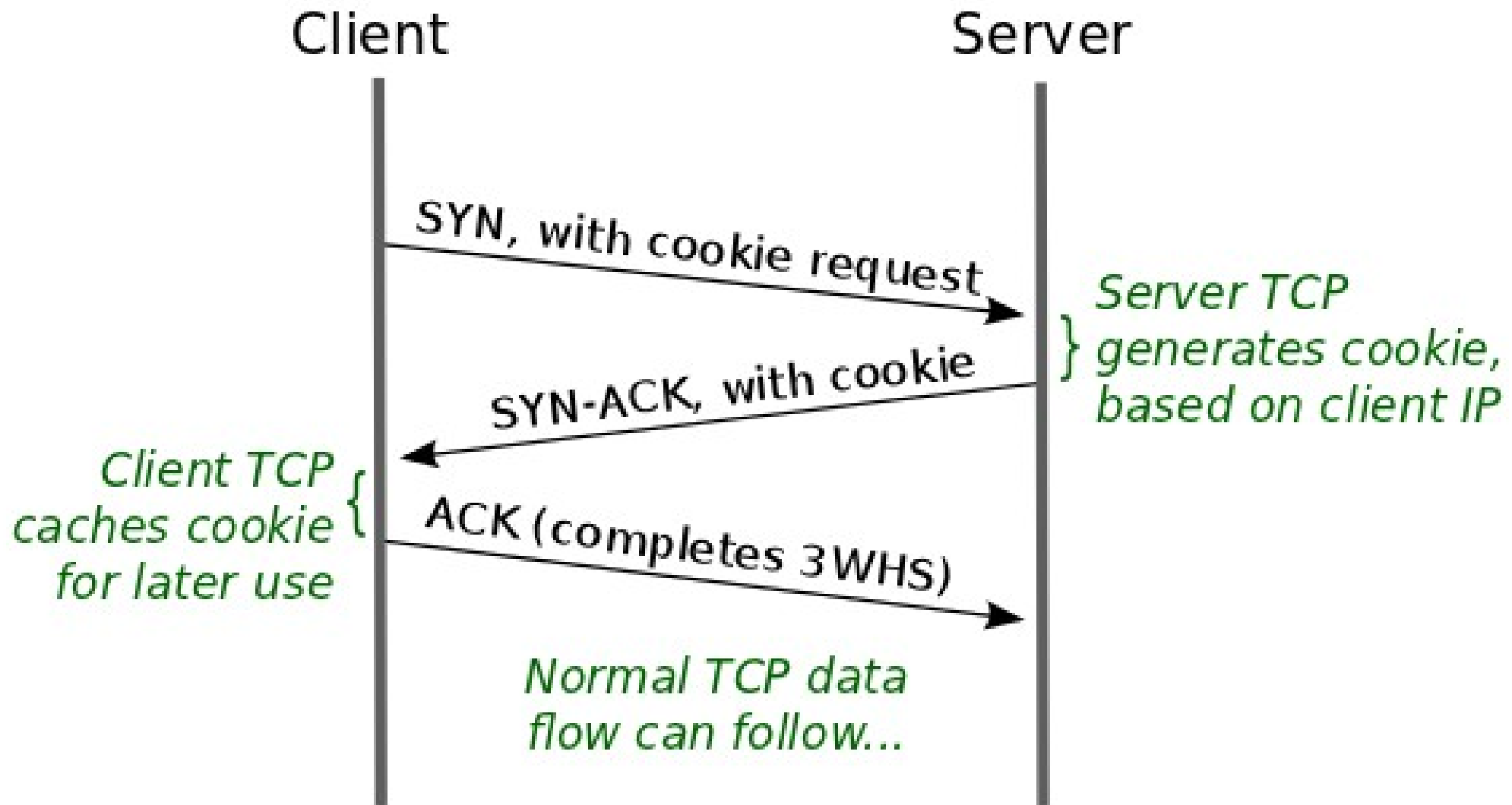# Transmission Datagram Protocol - TCP

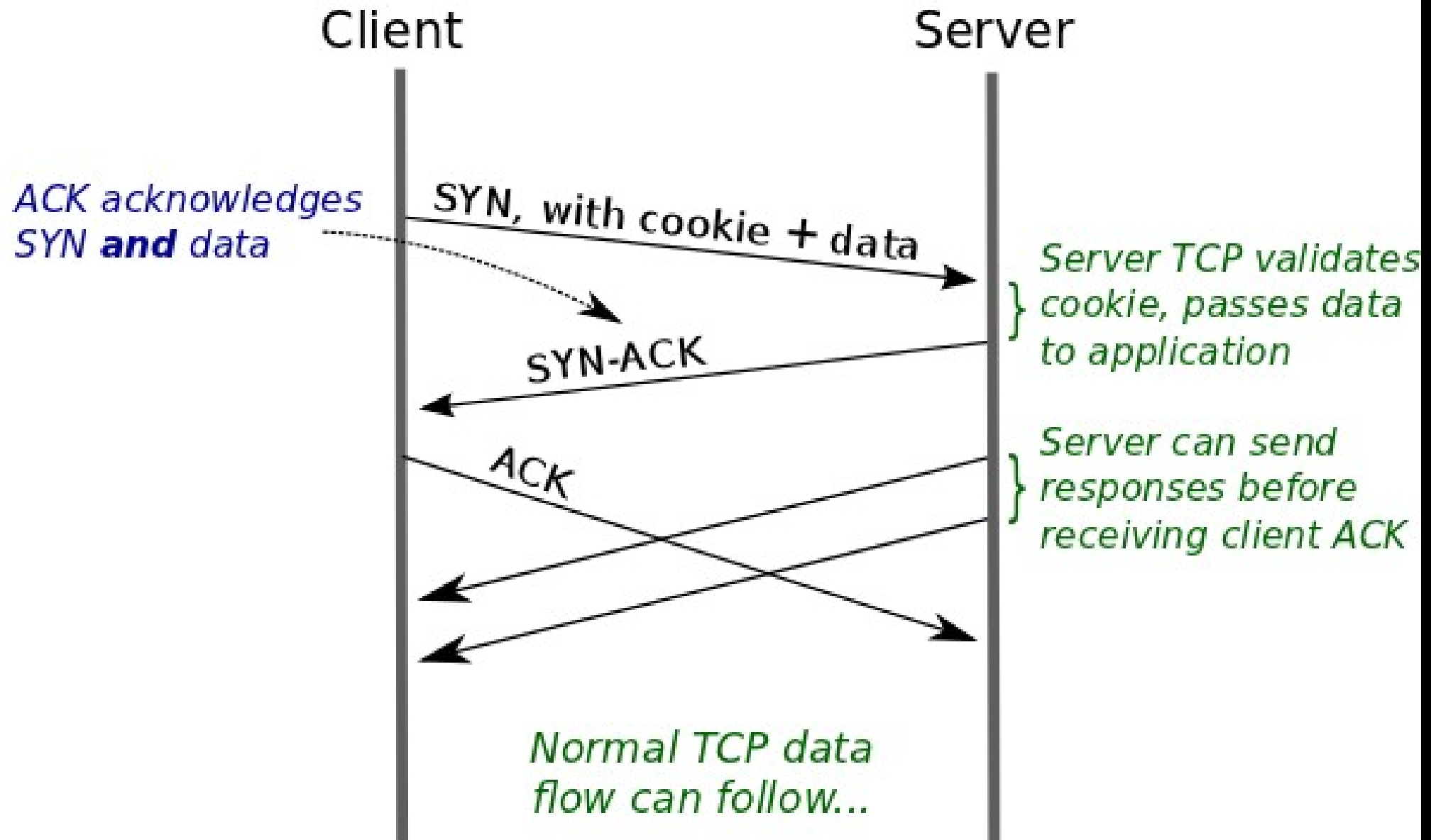## TCP Fast open

# Transmission Datagram Protocol - TCP

## TCP Fast open

# TCP Congestion

# TCP Congestion

- ➢ **Slow start**
- ➢ **Congestion avoidance**
- ➢ **Fast retransmit**
- ➢ **Fast Recovery**

# TCP Congestion – Slow start

➢ **The initial window size is initialized to one MSS**
➢ **Each time a packet is ACKed the congestion window i increased**
➢ **When the ssthresh is reached, the next phase starts**

# TCP Congestion –
# Congestion avoidance

➢ **In this phase window size is increased linearly until timeout occurs or duplicate ACK is received**

# TCP Congestion – Fast retransmit

➢ **If more then 3 ACKs are received for the same segmen[t] the sender has to send that particular segment even b[efore] its timer has expired**

# TCP Congestion – Fast Recovery

➢ **In this phase window size is decreased to ssthresh rat
then the smaller initial value and increase its size line**



http://histrory.visualland.net/tcp_fast_recovery.html

# TCP Congestion Avoidance - Problems

➢ Slow-start assumes that unacknowledged segments a
due to network congestion, which is usually NOT the c
in wireless networks, where dropped packets are mai
because of poor data link quality.

➢ The slow-start protocol performs badly for short-live
connections, because it actually slows down the
transmission of data.

➢ It is possible to trick the congestion avoidance algori
to think that the pipe is full and slow down all connect
originating from that machine.

# TCP Keepalive

➢ The keepalive packets are packets which contain no sent at regular interval to confirm that this connection alive

➢ Keepalive time is the duration between two keepalive transmissions in idle condition. TCP keepalive period required to be configurable and by default is set to no than 2 hours.

➢ Keepalive interval is the duration between two succe keepalive retransmissions, if acknowledgement to the previous keepalive transmission is not received. Usual around 75 seconds.

➢ Keepalive retry is the number of retransmissions to b sent out before declaring that remote end is not availa

# Datagram Congestion Control Protocol

➢ Basically DCCP is UDP with congestion control mecha

It features
➢ Unreliable flows of datagrams
➢ Reliable handshakes for connection setup and teardow
➢ Negotiation of a suitable congestion control mechanism
➢ Acknowledgment mechanisms communicating packet l
➢ Path Maximum Transmission Unit (PMTU) discovery

➢ RFC4340

# DCCP header (x = 1)

```
 0                   1                   2                   3
 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|          Source Port          |           Dest Port           |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|  Data Offset  | CCVal | CsCov |           Checksum            |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|     |       |X|           |                                   .
| Res | Type  |=|  Reserved |   Sequence Number (high bits)     .
|     |       |1|           |                                   .
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
.                  Sequence Number (low bits)                   |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

If X is 1 the Sequence Number field is 48 bits long

# DCCP header (x = 0)

```
 0                   1                   2                   3
 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|          Source Port          |           Dest Port           |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|  Data Offset  | CCVal | CsCov |            Checksum           |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|       |       |X|                                             |
| Res   | Type  |=|          Sequence Number (low bits)         |
|       |       |0|                                             |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

   If X is 0 the Sequence Number field is 24 bits long

Data Offset - the offset from the start of the packet's DCCP header to t
start of its application data area
CCVal - Defines the congestion control algorithm used
    CCVal = 2 - TCP like congestion avoidance
    CCVal = 3 - TCP friendly congestion avoidance
CsCov - Checksum Coverage determines the parts of the packet that are
covered by the Checksum field.
Checksum – DCCP header checksum
Type – DCCP packet type
X - Extended Sequence Numbers (may be 0 or 1)

# DCCP header

All currently defined packet types except DCCP-Request and DCCP-Data carry an Acknowledgment Number Subheader

When X=1, its format is:

```
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|              Reserved             |       Acknowledgment Number    .
|                                   |            (high bits)         .
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
.                 Acknowledgment Number (low bits)                |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```
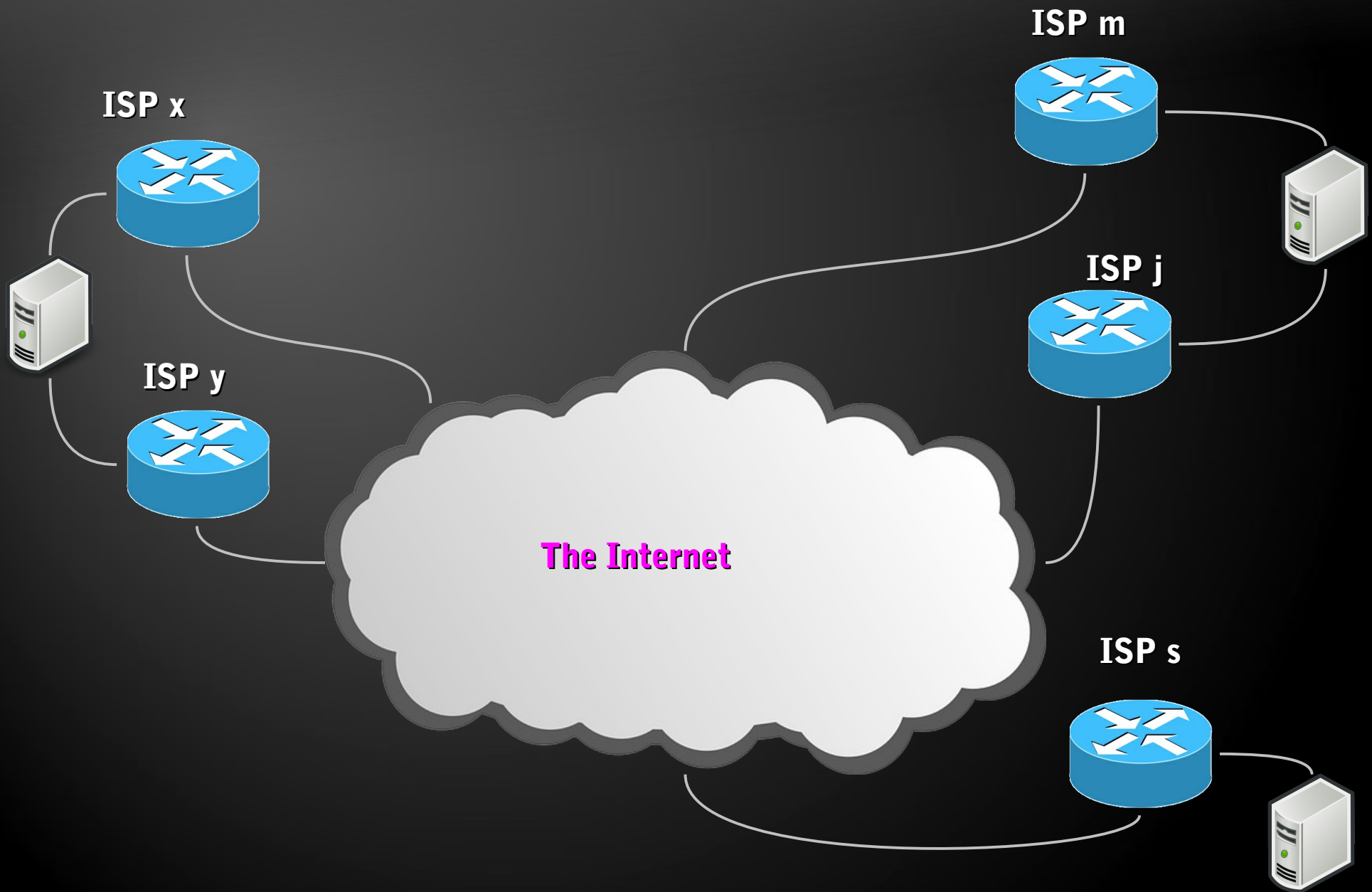
When X=0, only the low 24 bits of the Acknowledgment Number are transmitted, giving the Acknowledgment Number Subheader this format:

```
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|   Reserved    |        Acknowledgment Number (low bits)        |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

# DCCP packet types

```
Type    Meaning
----    -------
   0    DCCP-Request
   1    DCCP-Response
   2    DCCP-Data
   3    DCCP-Ack
   4    DCCP-DataAck
   5    DCCP-CloseReq
   6    DCCP-Close
   7    DCCP-Reset
   8    DCCP-Sync
   9    DCCP-SyncAck
10-15   Reserved
```
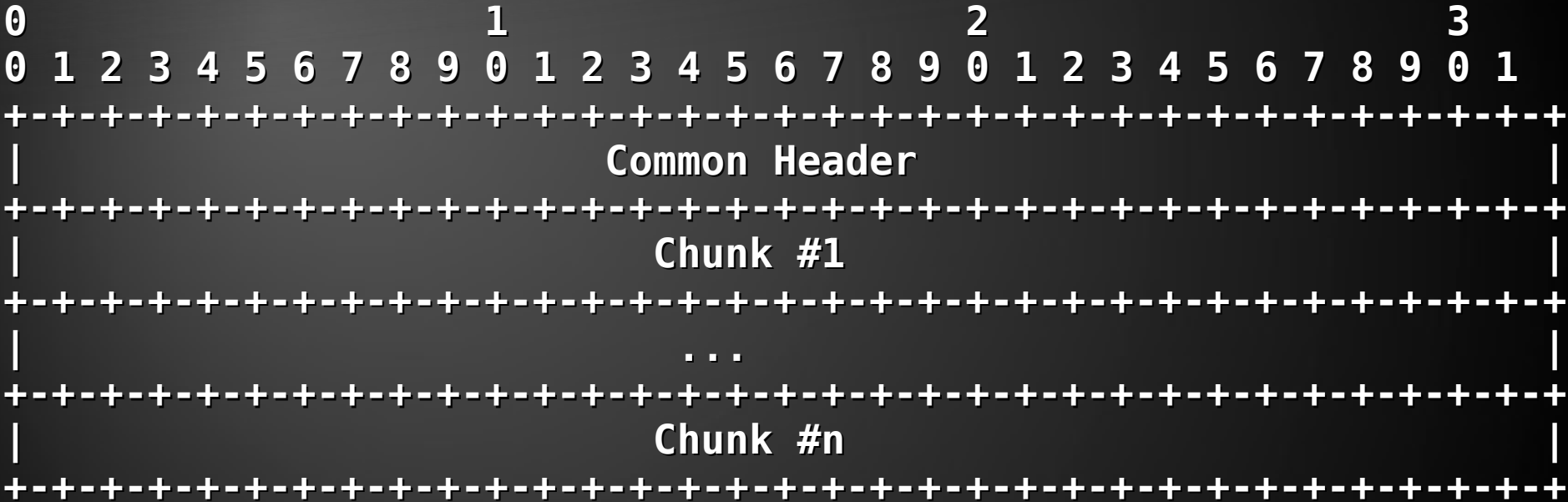
Multihoming

# Stream Control Transmission Protocol

➢ **Multihoming support in which one or both endpoints of a connection can consist of more than one IP address, enabling transparent fail-over between redundant network paths.**

➢ **Delivery of chunks within independent streams eliminate unnecessary head-of-line blocking**

➢ **Path selection and monitoring**

➢ **Validation and acknowledgment mechanisms protect against flooding attacks and provide notification of duplicated or missing data chunks.**

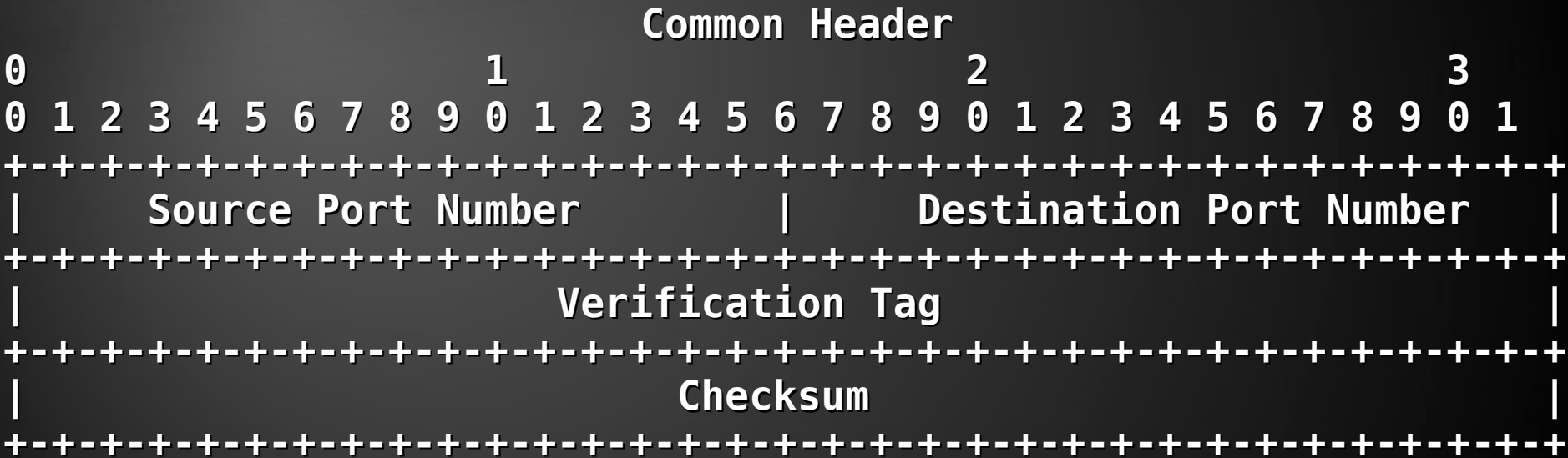➢ **Improved error detection suitable for Ethernet jumbo frames.**

# Stream Control Transmission Protocol

```
0                   1                   2                   3
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                         Common Header                         |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                          Chunk #1                             |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                            ...                                |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                          Chunk #n                             |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

## RFC4960

# Stream Control Transmission Protocol

```
                        Common Header
0                   1                   2                   3
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|     Source Port Number        |     Destination Port Number   |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                      Verification Tag                         |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                         Checksum                              |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

**RFC4960**

# IP & Domain allocation

- IANA – www.iana.org
  - Number resources
    - IP Addresses
    - Autonomous System (AS)
    - Protocol number assignments
  - Domain assignments
    - Root zone management
    - gTLD database
    - .int and .arpa domains
- IP registrars
  - ARIN, LACNIC, Africa, APNIC, RIPE

ARIN

RIPE NCC

AfriNIC

APNIC

LACNIC

# IP & Domain allocation

➢ Regional Internet Registrar(RIR)
➢ Local Internet Registrar(LIR)
➢ There are two types of IP addresses that can be reques
   ➢ Provider dependent
      ➢ These you get from your ISP
   ➢ Provider independent
      ➢ You get them from the local LIR or the regional RIR
      ➢ These allocations can not be smaller then /24 netwo
➢ Autonomous System (AS)
   ➢ Used for the BGP routing protocol
   ➢ Aggregated IP route announcements are made from
     to them
   ➢ The corner stone of the Internet routing
   ➢ Look at http://www.youtube.com/watch?v=oK-lgjJhC4

# Domain Name System - DNS

**Everything was 'hosts':**

```
127.0.0.1              localhost
192.168.0.174    store1
192.168.0.238    store2
192.168.0.244    store3
192.168.155.2    operations
192.168.155.149 zimbra0.siteground.com
193.107.36.190     sapport.bg www.sapport.bg
8.8.8.8              ns.google.com
89.25.120.31       google.com
89.25.120.24       www.google.com
```

**Linux: /etc/hosts**
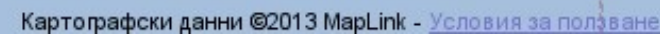**Windows: C:\Windows\System32\drivers\etc\hosts**

# Domain Name Space

"zone delegation"

**NS RR** ("resource record")
names the nameserver
authoritative for
delegated subzone

"delegated subzone"

When a system administrator
wants to let another administrator
manage a part of a zone, the first
administrator's nameserver **delegates**
part of the zone to another
nameserver.

**resource records**
associated with name

**zone** of authority,
managed by a **name server**

see also: RFC 1034 4.2:
How the database is divided into zones.

# DNS

- ➢ **Internet Corporation For Assigned Names and Numbe**
  - ➢ IANA is now part of it
  - ➢ Handles ccTLDs
  - ➢ Handles gTLDs
  - ➢ Handles the Root zone

- ➢ Country level domains
  - ➢ .bg, .co.za, .co.uk, .edu.us

- ➢ Top-level domains
  - ➢ .com, .net, .org, .edu, .gov, .mil
  - ➢ .biz, .name, .info

- ➢ **Instances of J and L root servers are hosted in Sofia**

# DNS

- ➢ **Internet Corporation For Assigned Names and Numbe**
  - ➢ **IANA is now part of it**
  - ➢ **Handles ccTLDs**
  - ➢ **Handles gTLDs**
  - ➢ **Handles the Root zone**
- ➢ **Internationalized domain name (IDN)**
  - ➢ **Domain names are encoded using** **Punycode**
  - ➢ **.ru = .рф**
  - ➢ **We are expecting soon .bg = .бг**
- ➢ **Country level domains**
  - ➢ **.bg, .co.za, .co.uk, .edu.us**

- ➢ **Top-level domains**
  - ➢ **.com, .net, .org, .edu, .gov, .mil**
  - ➢ **.biz, .name, .info**

# DNS

- ➤ **Name servers**
  - ➤ **Authoritative only**
  - ➤ **Recursive**
  - ➤ **Authoritative + recursive**
- ➤ **.in-addr.arpa**
- ➤ **.ip6.arpa**

Local NameServer

TLD Name Servers

INTERNET

NS1-4.GOOGLE.COM

www.google.com

# DNS - Resolving



**Default service port TCP/UDP: 53**

# DNS - Resolving

- **Forward resolving**
  - **Host/FQDN to IP**

- **Reverse resolving**
  - **IP to Host**

- **Reverse resolver delegation**
  - **RIR -> LIR -> Local ISP -> YOU**

# DNS

# DNS Resource records

| TYPE | value and meaning |
|------|-------------------|
| A | 1 a host address |
| NS | 2 an authoritative name server |
| CNAME | 5 the canonical name for an alias |
| SOA | 6 start of a zone of authority |
| WKS | 11 a well known service description |
| PTR | 12 a domain name pointer |
| HINFO | 13 host information |
| MINFO | 14 mailbox or mail list information |
| MX | 15 mail exchange |
| TXT | 16 text strings |
| AXFR | 252 A request for a transfer of an entire zone |

# DNS Resource records

```
kar-do.cc. 86400 IN SOA       ns1.ex1.com.
   mm.yuhu.biz.    (
      2013013106  ;Serial Number
      86400       ;refresh
      7200        ;retry
      3600000     ;expire
      86400       ;minimum
)
kar-do.cc.    IN      NS      ns1.ex1.com.
kar-do.cc.    IN      NS      ns2.ex1.com.
kar-do.cc.    IN      A       134.154.23.12
localhost     IN      A       127.0.0.1
kar-do.cc.    IN      MX      0  mail.kar-do.cc.
mail          IN      CNAME   mail.yuhu.biz.
www           IN      A       134.154.23.12
www           IN      A       134.142.65.81
kar-do.cc.    IN      TXT
         "v=spf1 +a +mx +ip4:134.154.23.12 ?all"
```

**RFC5321**

Send Mail Transport Protocol - SMTP

```
S: 220 smtp.example.com ESMTP Postfix
C: HELO relay.example.org
S: 250 Hello relay.example.org, I am glad to meet you
C: MAIL FROM:<bob@example.org>
S: 250 Ok
C: RCPT TO:<alice@example.com>
S: 250 Ok
C: RCPT TO:<theboss@example.com>
S: 250 Ok
C: DATA
S: 354 End data with <CR><LF>.<CR><LF>
C: From: "Bob Example" <bob@example.org>
C: To: "Alice Example" <alice@example.com>
C: Cc: theboss@example.com
C: Date: Tue, 15 January 2008 16:02:43 -0500
C: Subject: Test message
C:
C: Hello Alice.
C: This is a test message with 5 header fields and 4 lines in the message body.
C: Your friend,
C: Bob
C: .
S: 250 Ok: queued as 12345
C: QUIT
```

➢ **In this phase window size is increased linearly until timeout occurs or duplicate ACK is received**

Telerik Academy

# Questions?

Beer time

- **Operating Systems @ Telerik Academy**

- **http://telerikacademy.com/Courses/Courses/Details/35**

- **Telerik Software Academy**

  - **academy.telerik.com**

- **Telerik Academy @ Facebook**

  - **facebook.com/TelerikAcademy**

- **Telerik Software Academy Forums**

  - **forums.academy.telerik.com**