Lecture 12

# TCP/IP Protocol Suite II

**Martin Svoboda**
svoboda@ksi.mff.cuni.cz

24. 5. 2021

**Charles University**, Faculty of Mathematics and Physics

# Lecture Outline

**ICMPv4**
- Message structure and basic message types

**ARP**
- Translation of IP addresses to hardware addresses

**RARP**
- Translation of hardware addresses to IP addresses

**DHCP**
- Configuration of nodes including assignment of IP addresses

**IPv6**
- Packet structure, header fields, and fragmentation

**ICMPv6**

# ICMPv4 Protocol

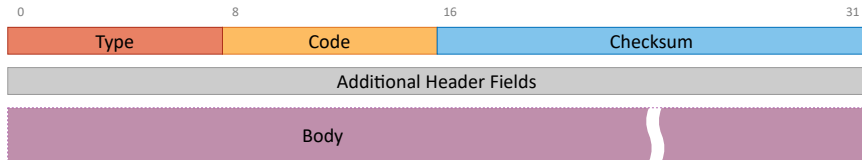**Internet Control Message Protocol** (**ICMP**)

- **Auxiliary <u>L3</u> protocol providing support to IPv4**
  - Allows to deal with **errors and non-standard situations**
    - Since IP alone is not capable of doing so
- **Types of messages**
  - **Error messages**
  - **Informational messages**: various queries, requests, replies, …
- **Scope of validity is not limited to just a single network**
  - And so ICMP messages need to be routed across networks
  - ICMP could thus work as yet another full-fledged L3 protocol
    - So that its messages would be inserted directly into L2 frames
    - But that would mean ICMP itself would need to be routable
  - And so instead, **ICMP messages are inserted into IP datagrams**
    - Which kind of (again) contradicts principles of layered models

# ICMPv4 Messages

**Message structure**

- **Header** (64 bits)
  - Always fixed length, though **partially variable header fields**
- **Body**
  - May entirely be missing
  - Fixed or variable length
    - This length does not need to be explicitly remembered
    - Since it is derivable from the length of the entire IP datagram

| 0 | 8 | 16 | 31 |
|---|---|---|---|
| Type | Code | Checksum | |
| Additional Header Fields | | | |
| Body | | | |

# ICMPv4 Messages

**Header fields**

- **Type** (8 bits)
  - **Main type** of a given ICMP message
- **Code** (8 bits)
  - Code describing a **particular subtype** of a given message type
- **Checksum** (16 bits)
  - **Checksum of the <u>whole</u> message**, not just its header
    - The same calculation mechanism as in case of IP itself is used
- **Additional header fields** (32 bits)
  - Depend on a particular message type
    - **Often unused** (but always present)

# ICMPv4 Messages

Message **types and codes**

- Maintained by **IANA**
  - https://www.iana.org/assignments/icmp-parameters/
  - Current state (as of May 2021)
    - Almost $\approx$ 45 out of 256 types are used or reserved
    - Many of which are deprecated, though

Message **body**

- **Depends on a particular message type**
  - May entirely be missing as already outlined
- Often contains **beginning of the original IP datagram**
  - It means IP datagram which caused a given error message
  - In particular, its **full header and the first 64 bits of its body**
    - So that **source and destination** L4 **ports** are available, too
    - Which helps in correct recognition of the original datagram

# Message Examples

**Destination Unreachable** (Type 3)

- Error messages sent when datagrams needed to be discarded
  - Since their further processing was not possible
    - Because of various particular reasons…
- **Network Unreachable** (Code 0)
  - Sent by routers on the way when…
  - **Network of the intended destination node is not reachable**
    - I.e., its distance in the routing table equals to infinity
    - Or the destination network is unknown at all
- **Host Unreachable** (Code 1)
  - Sent by inbound routers in the destination network when…
  - **Intended destination node is not reachable**
    - In case such a detection is possible at all

# Message Examples

**Destination Unreachable** (Type 3) (cont'd)

- **Protocol Unreachable** (Code 2)
  - Sent by recipient nodes when…
  - **Designated transport protocol is not supported**
    - More precisely, payload type in **protocol field** is not supported
- **Port Unreachable** (Code 3)
  - Sent by recipient nodes when…
  - **Specified destination** L4 **port is invalid or not available**
    - I.e., it is not possible to perform datagram demultiplexing
- **Fragmentation Needed** (Code 4)
  - Sent by routers on the way when…
  - **Fragmentation is needed but was explicitly prohibited**
    - I.e, when the **Don't Fragment Flag** was enabled
- …

# Message Examples

**Time Exceeded** (Type 11)

- **Time to Live Exceeded in Transit** (Code 0)
  - Sent by routers on the way when…
  - **Datagram TTL was exceeded**
    - TTL value dropped to zero

- **Fragment Reassembly Time Exceeded** (Code 1)
  - Sent by recipient nodes when…
  - **Reassembling of a fragmented datagram was not possible**
    - Since not all fragments were received within a given timeout

# Message Examples

**Echo Request** (Type 8, Code 0) and **Echo Reply** (Type 0, Code 0)

- **Allow for testing of reachability of nodes**
  - Sender sends an **Echo Request** message
  - Recipient responses with an **Echo Reply** message
    - At least should respond (it is compulsory)
    - But nowadays often does not because of security reasons
- Additional **header fields** are needed
  - **Allow to match corresponding pairs of requests and replies**
    - Since multiple requests may be sent in a row
  - **Identifier** and **Sequence Number**
    - **Unique identification and serial number given by a sender**
    - Particular implementation varies across individual systems
  - Both fields are preserved in replies, including message data
- Used by `traceroute` or `ping` utilities

# Ping Tool

`ping` (backronym **Packet InterNet Groper**)

- **Diagnostic tool** allowing for testing of **reachability of nodes**
  - Together with **measured round-trip delivery times**
  - As well as basic **datagram loss statistics**

Basic principle

- **ICMP Echo Request** is sent to the intended destination
  - Actually a whole **series of requests** is sent, one by one
    - In order to make the detection and measurement more precise
- When **ICMP Echo Reply** response is received
  - Intended destination is reachable
- When **no response is received within a timeout**
  - Destination is considered as unreachable
    - Though it might just be unwilling to respond

# Message Examples

**Source Quench** (Type 4, Code 0) (deprecated)

- **Feedback technique for Congestion Control and Flow Control**
  - When **routers or end nodes** reached their **capacity limits**
    - In terms of available buffer sizes, transmission capacity, …
  - Or even better, when they are just **approaching such limits**
- Principle
  - **Source nodes potentially causing the issues are informed**
    - I.e., respective sender or senders
  - So that **discarding of datagrams** is attempted to be avoided
- Observations
  - It is not possible to specify the **extent of congestion** problems
    - And so it is also not clear how the source nodes should react
  - Nor it is possible to revoke the original announcement

# Message Types

Other interesting message types

- **Redirect** (Type 5)
  - **Signals incorrect or not optimal routing**
    - At the level of particular nodes or whole networks
- **Router Advertisement** (Type 9)
  - **Allows routers to reveal their existence** to end nodes
    - Using multicast to 224.0.0.1 (or broadcast when unavailable)
- **Router Solicitation** (Type 10)
  - **Allows end nodes to search for available routers**
    - Using multicast to 224.0.0.2 (or broadcast when unavailable)
- **Parameter Problem** (Type 12)
  - Messages representing **generic so far not covered problems**
- …

# ARP

**Address Resolution Protocol** (**ARP**)

- Allows for **translation of IP addresses to hardware addresses**
  - For the purpose of **direct delivery** at L3
    - I.e., delivery within the context of a given network
- In particular…
  - We are about to **locally send an** L3 **IP datagram**
    - To the final recipient or just the first / next router on the way
    - **Intended recipient is expressed in terms of its** L3 **IP address**
    - I.e., in both the cases, we are provided with this address
  - Delivery of the datagram is, however, **executed using** L2
    - And L2 is only capable of working with L2 hardware addresses
    - Unfortunately, we only have the intended L3 IP address…
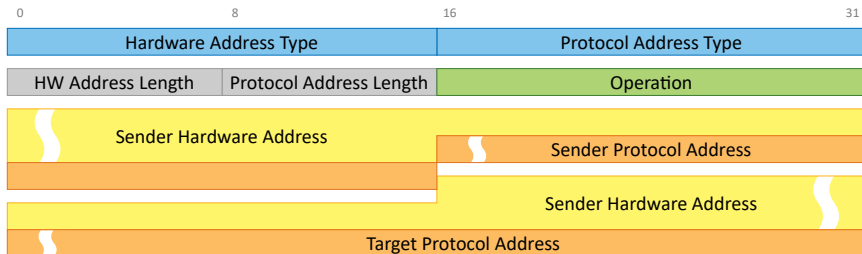  - And so the corresponding L2 **address needs to be discovered**

# ARP

**Address Resolution Protocol** (cont'd)

- Basic idea
    - **Sender creates an ARP Request message**
        - Includes the queried IP address in this message
        - And **sends it to the whole network using** L2 **broadcast**
    - **Corresponding target node captures this request**
        - By matching its IP address with the queried one
        - **Creates an ARP Reply message with its hardware address**
        - And **sends it back via an ordinary** L2 **unicast**
- Observations
    - **Different technologies and addresses are used at** L2
    - As well as **different protocols and addresses are used at** L3
- ⇒ it would be nice to handle them all in a unified way
    - ARP really is capable of such **universal applicability**

# ARP Messages

**Message structure**

- All fields are compulsory, overall length is variable
  - Since individual **address fields have variable lengths**
- **ARP messages are encapsulated to L2 frames**
  - And so ARP as such belongs to the L3 network layer

| 0 | 8 | 16 | 31 |
|---|---|---|---|
| Hardware Address Type | | Protocol Address Type | |
| HW Address Length | Protocol Address Length | Operation | |
| Sender Hardware Address | | | |
| | | Sender Protocol Address | |
| | | Sender Hardware Address | |
| Target Protocol Address | | | |

# Message Structure

**Hardware Address Type** (16 bits) and **Length** (8 bits)

- Describe the **type and length of** L2 **addresses**
  - I.e., identify a particular L2 technology
  - As well as length of its addresses in a number of bytes
- Types are maintained by **IANA**
  - https://www.iana.org/assignments/arp-parameters/
- Examples
  - **Ethernet** 10 Mb/s (type 1, length 6 bytes)
  - IEEE 802 Networks with **EUI-48** (type 6, length 6 bytes)
  - **EUI-64** (type 27, length 8 bytes)
  - **HDLC** (type 17, length 1 byte or more)
  - …

# Message Structure

**Protocol Address Type** (16 bits) and **Length** (8 bits)

- Describe the **type and length of** `L3` **addresses**
  - I.e., identify a particular `L3` protocol
  - And similarly length of its addresses
- Types are (primarily) maintained by **IEEE RA**
  - **EtherTypes** were recycled for this purpose
  - http://standards-oui.ieee.org/ethertype/eth.txt
- Examples
  - **IPv4** (type `0x0800`, length 4 bytes)
  - **IPv6** (type `0x86DD`, length 16 bytes)

# Message Structure

**Operation** (16 bits)
- Allows to distinguish individual **ARP operations**
- Codes are maintained by **IANA**
  - https://www.iana.org/assignments/arp-parameters/
  - **Request** (1), **Reply** (2), …

**Sender Hardware Address** and **Sender Protocol Address**
- L2 and L3 addresses of the sender
  - I.e., node sending a given request or reply

**Target Hardware Address** and **Target Protocol Address**
- L2 and L3 addresses of the indented recipient

# Resolution Process

**ARP Request message** (operation code 1)

- **Type and length fields** are set according to the situation
- **Address fields** are filled in as follows…
  - Sender hardware and protocol addresses
    - Both are set according to the sender
  - **Target hardware address** is left undefined
  - **Target protocol address** is set to the **queried IP address**

**ARP Reply message** (operation code 2)

- Reply can be constructed directly from the received request
  - **Operation code** is changed from request to reply
  - **Source and target addresses are mutually swapped**
  - **Source HW address** is then set to the **resolved HW address**

# Resolution Process

**Resolution** process has **significant overhead**

- Not just because **broadcast is required**
- It would also be inefficient to repeat requests over and over
  - And so **discovered mappings are cached**

## ARP Cache

- **Table** with resolved **IP and hardware address bindings**
- **Static** records (e.g.: 192.168.1.255 → FF-FF-FF-FF-FF-FF)
- **Dynamic** records
  - Must be **periodically forgotten**
    - So that changes within the network can be reflected
    - Timeout can be 1 minute for end nodes, hours for routers
  - As well as **refreshed** to restrict new unnecessary queries
    - With the aim of optimizing the whole process even more

# Resolution Process

**Resolution** steps

- **ARP Cache table is first consulted**
  - When the required **binding already exists**
    - It is simply fetched and the whole process ends
- **Otherwise** an **ARP Request message** is constructed
  - And sent to the whole network using L2 broadcast
- **Each and every node captures the request message**
  - And exploits the received information to update its cache
    - I.e., **adds a new binding** or **refreshes an already existing one**
- **Intended target node** (if any) furthermore…
  - Creates an **ARP Reply message**
  - And sends it to the original node using L2 unicast

# Reverse ARP

**Reverse Address Resolution Protocol** (**RARP**)

- Allows for **translation of hardware addresses to IP addresses**
  - In terms of **assignment of an IP address to a given node**

**Assignment** process

- **RARP Request message** is constructed (**operation** code 3)
  - Message format remains the same as in the traditional ARP
    - Both **hardware address fields** are set to the known HW address
    - Both protocol address fields are left unused
  - **Request is then sent to the whole network via** L2 **broadcast**
- **RARP Server** captures this request
  - I.e., special host configured to serve such requests
  - **RARP Reply message** is constructed (**operation** code 4)
  - Reply is then sent back to the original node via L2 unicast

# Reverse ARP

**Drawbacks**

- Very old and simple solution
- **RARP operates at** `L3`
  - Encapsulates its messages directly into `L2` frames
  - And so **RARP server must be available in each network**
    - Since RARP messages cannot cross network boundaries
- Whole approach **cannot work without** `L2` **broadcast**
- **Only fixed manually defined bindings are supported**
  - This is not sufficient from today's perspective
- **Additional information cannot be passed at all**
  - I.e., only IP address itself can be assigned
  - And not other (nowadays) essential information
    - E.g., netmask / CIDR prefix, router IP address, …

# DHCP

**Dynamic Host Configuration Protocol** (**DHCP**)

- Newer solution dealing with the identified drawbacks
  - Based on **BOOTP** (**Bootstrap Protocol**)
    - Allow for IP address assignment, but its primary motivation was related to providing boot images to diskless workstations
- Advantages
  - One DHCP server can serve **multiple networks**
  - **Dynamic assignments** of addresses is possible
  - **Operates at** L7, uses UDP datagrams at ports 67 and 68
  - Allows for interchange of **additional information**
    - Nemask, routers, DNS servers, time zone, time servers, …
- Three **allocation strategies** are provided
  - **Manual**, **Automatic** and **Dynamic**

# Allocation Strategies

**Manual Allocation** (also **Static Allocation**)

- Requesting client acquires a **predefined IP address**
  - Based on fixed **HW address → IP address bindings**
  - Created manually by the **network administrator** in advance
- As a consequence, **allocated address is always the same**
  - Which is suitable for network printers or similar devices

**Automatic Allocation**

- Requesting client acquires an **arbitrary IP address**
  - Chosen by the DHCP server itself from a given **address pool**
- Allocation is understood as **permanent**
  - It means the **binding is remembered the first time it is created**
  - So that the next time **the same address** can be granted again

# Allocation Strategies

## Dynamic Allocation

- Requesting client acquires an **arbitrary IP address**
  - Again chosen by the DHCP server from a given **address pool**
    - Of course, only currently unused addresses can be considered
- However, allocation is **<u>temporary</u>** only in this case…
  - Based on a **concept of lease**
    - I.e., only for a **limited period of time**
    - Which is specified at the moment of the allocation
- As a consequence…
  - **Different address may be provided each time!**
    - As well as one address can gradually be used by different nodes
- This has a **fundamental impact on IP address management**
  - Though it may not be apparent at first sight…

# Lease Concept

**Traditional approach**

- Once assigned, **nodes have their addresses permanently**
  - In a sense of being their **owners or holders**
  - And so they can use them as long as they want to
    - I.e., for any length of time without any limitation
- This approach is **simple** on one hand
- But, unfortunately, **not flexible enough** on the other
  - Since **end nodes often roam from one network to another**
    - And so the whole traditional concept is no longer suitable for contemporary networks

# Lease Concept

Newly introduced **concept of lease**

- Nodes act as **DHCP clients** in a sense of **temporary** lessees
  - **They must proactively take care of their IP addresses**
    - They are expected to perform **various tasks**
    - And so make transitions through **various states**

**Period of lease**

- Appropriate length depends on a particular situation
  - **Shorter periods**
    - **Higher flexibility**, **lower stability**, **higher overhead**
  - **Longer periods**… on the contrary
- In practice…
  - Hours, days, weeks, months, …

# Client Actions

**Allocation**

- **Client does not yet have an IP address and asks for a lease**

**Reallocation**

- **Client does have an IP address and asks for its confirmation**
  - I.e., lease of the current address is still valid
    - And so there is actually no reason for such a request
- However, it is **voluntarily willing to accept a new address**
  - Since change at this moment would not cause any obstacles
  - E.g., because this node…
    - Has **just rebooted** or was **turned off for some time**

# Client Actions

**Renewal**

- **Client has an IP address but its lease is approaching its end**
  - And so **extension of lease period is requested**
    - If granted, a new lease with the same address is in fact started
- First renewal attempt is initiated in 50% of lease time

**Rebinding**

- **Client asks a different server for the currently leased address**
  - Suitable when the original server became unavailable
    - I.e., when standard renewal could not be finished successfully
- First rebinding attempt is initiated in 87.5% of lease time

**Release**

- **Client returns its IP address before its lease expired**

# IPv6 Protocol

**Internet Protocol version 6** (**IPv6**)
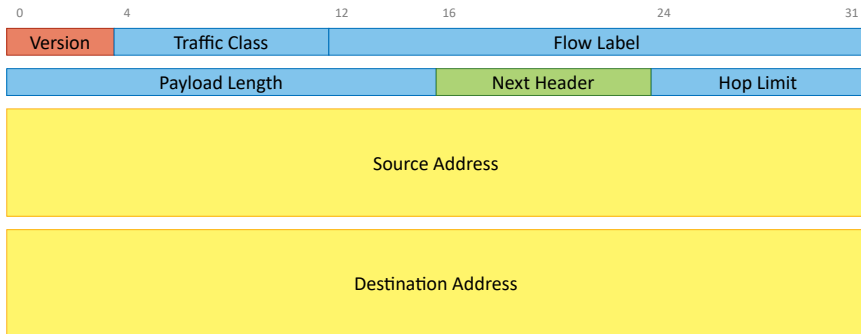
- **Differences** with respect to IPv4
    - **Larger IPv6 addresses**
        - 128 bits instead of just 32 bits
        - Together with 3 levels of routing (site / subnet / interface)
    - Simpler packet format
        - **Lower number of header fields**
        - Meaning and / or names of certain fields were changed
        - Some were removed entirely, e.g., **header checksum**
    - Concept of **extension headers**
        - Instead of IPv4 options
    - Different approach to **fragmentation**
    - Integrated **QoS support**
    - …

# IPv6 Packets

**Packet structure**
- **Header chain**
  - **Main header** (40 bytes) and optional **extension headers**
- Optional **body**

| 0 | 4 | 12 | 16 | 24 | 31 |
|---|---|---|---|---|---|

| Version | Traffic Class | Flow Label | | | |
|---|---|---|---|---|---|
| Payload Length | | Next Header | | Hop Limit | |
| Source Address | | | | | |
| Destination Address | | | | | |

# Main Header Fields

**Version** (4 bits)
- Fixed value 6

**Payload Length** (16 bits)
- Overall size of **payload and extension headers** (if any)
  - Main header is not included

**Hop Limit** (8 bits)
- Analogy to IPv4 **Time to Live** field

**Traffic Class** (8 bits)
- Analogy to IPv4 **Type of Service** field
  - I.e., used for the purpose of **Differentiated Services**

# Main Header Fields

**Flow Label** (20 bits)

- Allows to identify a particular **flow** = **group of related packets**
  - With the aim of treating them all in a similar way
    - E.g., with respect to QoS or other purposes
- In fact, **(Source Address, Flow Label)** forms the full identifier
  - Which allows to recognize such flows even at L3
- In IPv4, **transport connection** identification would be needed
  - I.e., tuple (sender $IP_1$:$port_1$, protocol, recipient $IP_2$:$port_2$)
    - Which is less convenient when compared to IPv6
    - Since L4 fields would need to be accessed in the payload

# Extension Headers

**Packet structure**

- Whole packet is composed from a **chain** of…
    - **Compulsory main header**
    - **Arbitrary number of extension headers** (0 or more)
        - Each should **only be used at most once** (exception exists)
        - They should be used in a specific **recommended order**
        - So that processing of IPv6 packets by routers is simplified
    - **Optional body**
- Enumerated blocks are put into the packet one after another
    - **Each header contains the Next Header field**
        - Which allows to mutually chain the individual blocks
        - I.e., describe what the next block is supposed to contain

# Extension Headers

**Next Header** (8 bits)

- Determines the **type of the next block** in a chain
  - I.e., type of the next **extension header** or **body payload**
  - Assuming that body (if any) must be placed at the very end
- Types (**Protocol Numbers**) are maintained by **IANA**
  - https://www.iana.org/assignments/protocol-numbers/
  - For simplicity, codes correspond to IPv4 analogies
- **Examples**
  - Extension headers: **IPv6 Fragmentation** (44), …
  - Payload protocols: **UDP** (17), **TCP** (6), **ICMPv6** (58), …
  - Special type: **IPv6 No Next Header** (59)
    - Suggests that noting follows
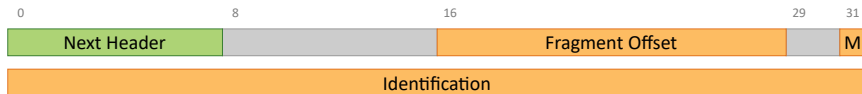    - And even if something does follow, it must be ignored

# IPv6 Fragmentation

**Fragmentation** in IPv6

- Differences with respect to IPv4
  - **Only source nodes can perform fragmentation** (never routers)
    - So that they can focus on their primary objective
    - Excessive packets are then automatically discarded
  - All information is stored within the **Fragmentation Header**
    - I.e., it is only used when fragmentation really took place
  - **Guaranteed minimal non-fragmenting MTU** is 1280 B
    - Compared to just 68 B / 576 B in case of IPv4
  - **IPv6 Path MTU Discovery**
    - Basically the same idea as in IPv4, though differences exist…
    - **ICMPv6 Packet Too Big** message is received instead
    - Includes value of the particular **MTU that caused the problem**

# IPv6 Fragmentation
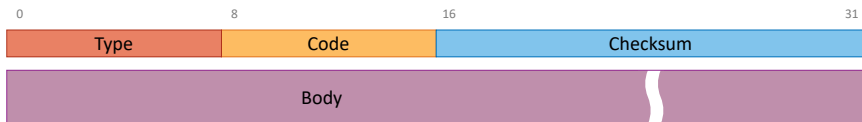
**Fragmentation header** (type 44)

- **Next Header** (8 bits)
- **Fragment Offset** (13 bits)
  - **Relative to the end of the last non-fragmented header**
    - I.e., certain extension headers are fragmented, others not
    - Fragmentation header is then put in between these two groups
- **More Fragments Flag** (1 bit)
- **Identification** (32 bits)
  - The same principle as in IPv4, only larger range

| 0 | 8 | 16 | 29 | 31 |
|---|---|---|---|---|
| Next Header | | Fragment Offset | | M |
| Identification | | | | |

# ICMPv6

**Internet Control Message Protocol version 6** (**ICMPv6**) (type 58)

- **Analogy to ICMPv4 for IPv4**
  - Basic principles are the same, though differences exist…
- **Longer part of the original IP packet** is preserved in body
  - As much as can be included not to exceed packet size 1280 B
    - So that fragmentation is avoided
- **Checksum calculation** also involves ICMP **pseudo-header**
  - With IPv6 source / destination addresses and other fields
- Slightly different generic **message structure**
- Different **types of particular ICMP messages**

| 0 | 8 | 16 | 31 |
|---|---|---|---|
| Type | Code | Checksum | |
| Body | | | |

# Lecture Conclusion

**ICMPv4**
- Destination Unreachable, Time Exceeded, …

**ARP** and **RARP**
- Translation of IP addresses to hardware addresses or vice versa

**DHCP**
- Manual, automatic, and dynamic **allocation strategies**
- Concept of **lease**, client actions

**IPv6**
- Main header, extension headers, body
- Fragmentation

**ICMPv6**