



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Introduction to RAW-sockets

Jens Heuschkel, Tobias Hofmann, Thorsten Hollstein, Joel Kuepper

16.05.2017

Technical Report No. TUD-CS-2017-0111
Technische Universität Darmstadt

Telecooperation Report No. TR-19,
The Technical Reports Series of the TK Research Division, TU Darmstadt
ISSN 1864-0516

<http://www.tk.informatik.tu-darmstadt.de/de/publications/>

Introduction to RAW-sockets

by

Heuschkel, Jens

Hofmann, Tobias

Hollstein, Thorsten

Kuepper, Joel

May 17, 2017

Abstract

This document is intended to give an introduction into the programming with RAW-sockets and the related PACKET-sockets. RAW-sockets are an additional type of Internet socket available in addition to the well known DATAGRAM- and STREAM-sockets. They do allow the user to see and manipulate the information used for transmitting the data instead of hiding these details, like it is the case with the usually used STREAM- or DATAGRAM sockets. To give the reader an introduction into the subject we will first give an overview about the different APIs provided by Windows, Linux and Unix (FreeBSD, Mac OS X) and additional libraries that can be used OS-independent. In the next section we show general problems that have to be addressed by the programmer when working with RAW-sockets. We will then provide an introduction into the steps necessary to use the APIs or libraries, which functionality the different concepts provide to the programmer and what they provide to simplify using RAW and PACKET-sockets. This section includes examples of how to use the different functions provided by the APIs. Finally in the additional material we will give some complete examples that show the concepts and can be used as a basis to write own programs. The examples are programmed in C++ and we assume that the reader has basic programming skills and networking knowledge to be able to understand the listings and content of this document.

Contents

1	Introduction	8
1.1	RAW-sockets	8
1.2	PACKET-sockets and Data Link Layer APIs	9
2	Implementations for different Operating Systems	10
2.1	Windows	10
2.1.1	Winsock-API	10
2.1.2	wincap	10
2.2	Linux	10
2.2.1	RAW-sockets	10
2.2.2	PACKET-sockets	11
2.3	Unix (FreeBSD, Mac OS X)	11
2.3.1	RAW-sockets	11
2.3.2	Berkeley Packet Filter (BPF)	12
2.4	OS independent	12
2.4.1	pcap	12
2.4.2	libnet	12
3	Programming with the APIs	13
3.1	General	13
3.1.1	Byte Order	13
3.1.2	Checksum	13
3.1.3	Type-Casting	14
3.1.4	Linux	14
3.1.4	Header-Positions	14
3.2	RAW-sockets	15
3.2.1	socket()	15
3.2.2	Unix	16
3.2.2	setsockopt()	16
3.2.2	Unix	17
3.2.3	getsockopt()	17
3.2.3	Unix	18
3.2.4	bind()	18
3.2.4	Unix	18
3.2.5	getsockname()	19
3.2.5	Unix	19
3.2.6	connect()	19
3.2.6	Unix	19
3.2.7	Read	19
3.2.7	read()	20
3.2.7	recv()	20
3.2.7	recvfrom()	20
3.2.7	flags and errors for the recv()- and recvfrom()-functions	21
3.2.7	Unix	22
3.2.8	Write	22
3.2.8	write()	22
3.2.8	send()	22
3.2.8	sendto()	23
3.2.8	Flags and errors for the send() and sendto() functions	23
3.2.8	Unix	24
3.2.9	close()	24
3.2.9	Unix	25
3.2.10	inet_ntop()	25
3.2.10	Unix	25
3.2.11	inet_pton()	25
3.2.11	Unix	26
3.2.12	Data Types	26
3.2.13	Layer 4	28

	Unix	28
	Read	28
	Write	30
3.2.14	Layer 3	30
	Read	32
	Write	32
3.2.15	Layer 2 - PACKET-sockets	36
	Promiscuous Mode	36
	MAC-Address	36
	Read	36
	Write	38
3.3	Berkeley Packet Filter (BPF)	40
3.3.1	BPF Header	41
3.3.2	Buffer Modes	41
3.3.3	IOCTLS	41
3.3.4	SYSCTL Variables	42
3.3.5	Filter Maschine	42
3.3.6	Read	42
3.3.7	Write	44
3.4	Winsock-API	44
3.4.1	Preparations and Usage	46
4	Programming with the libraries	47
4.1	Libnet	47
4.1.1	Preparations	47
4.1.2	Process of packet creation	47
	Library initialization	47
	Packet building	49
	Packet write	49
	Packet destruction	49
4.1.3	Functions	51
	Library initialization	51
	Build Functions	51
	Write	54
	Destruction	54
	Other functions	55
4.2	libpcap	58
4.2.1	Preparation	58
4.2.2	Process of packet capturing	58
	Interface Selection	58
	Initialize libpcap session	58
	Define Filters	58
	Initiate Sniffing process	59
	Identification of the header size	59
	Casting	60
	Byte order	60
	Close libpcap session	60
4.2.3	Functions	60
	Interface Selection	60
	Initialize libpcap session	62
	Define Filters	62
	Initiate Sniffing process	63
	Close Session	65
	Other functions	65
4.3	libpcap in Windows	65
4.3.1	Installation and Preparation	66
4.3.2	Process of packet capturing	66
4.3.3	Functions	66
5	Literature	67

A	Appendix: Listings	69
A.1	rfc1071 checksum cpp	69
A.2	win socket cpp	70
A.3	libnet tcp c	72
A.4	sendpack c	75
A.5	dump c	76
B	Appendix: Tables	80
B.1	Protocol Types of <netinet/in.h>	80
B.2	Linux Protocol Types defined in linux if_ether.h	81
B.3	Socket level options for setsockopt()	82
B.4	IP level options for setsockopt()	83
B.5	Errno flags for connect()	84
B.6	ioctl() flags defined in bpf.h	85

List of Figures

1	Socket provided by the operating system	8
2	Overview over the layers and access possibilities	9
3	Structure of a RAW-socket layer 4 Read Operation	29
4	Structure of a RAW-socket layer 4 Write Operation	31
5	Structure of a RAW-socket layer 3 Read Operation	33
6	Structure of a RAW-socket layer 3 Write Operation	34
7	Structure of a PACKET-socket layer 2 Read Operation	37
8	Structure of a PACKET-socket layer 2 Write Operation	39
9	Structure of a BPF device layer 2 Read Operation	43
10	Structure of a BPF device layer 2 Write Operation	45
11	Overview of the packet construction in libnet	48
12	Overview of the packet construction in libnet	50

List of Tables

1	Byte-Order Transformation Functions [8]	13
2	C-Header-Files for Network Headers [8]	14
3	Address family constants provided in <code><sys/socket.h></code> [8]	15
4	Errno flags for <code>socket()</code> as defined in <code><errno.h></code> [8]	15
5	The socket types are defined in <code><sys/socket.h></code> [8]	16
6	Errno flags for <code>setsockopt()</code> as defined in <code><errno.h></code> [8]	17
7	Errno flags for <code>getsockopt()</code> as defined in <code><errno.h></code> [8]	18
8	Errno flags for <code>bind()</code> as defined in <code><errno.h></code> [8]	18
9	Errno flags for <code>UNIX-bind()</code> as defined in <code><errno.h></code> [8]	18
10	Errno flags for <code>getsocketname()</code> as defined in <code><errno.h></code> [8]	19
11	Errno flags for <code>read()</code> as defined in <code><errno.h></code> [8]	20
12	Flags for <code>recv()</code> and <code>recvfrom()</code> defined in <code><sys/socket.h></code> [8]	21
13	Errno flags for <code>recv()</code> and <code>recvfrom()</code> as defined in <code><errno.h></code> [8]	21
14	Errno flags for <code>write()</code> as defined in <code><errno.h></code> [8]	22
15	Flags for <code>send()</code> and <code>sendto()</code> as defined in <code><errno.h></code> [8]	24
16	Errno flags for <code>send()</code> and <code>sendto()</code> as defined in <code><errno.h></code> [8]	24
17	Errno flags for <code>close()</code> as defined in <code><errno.h></code> [8]	25
18	Errno flags for <code>ntop()</code> as defined in <code><errno.h></code> [8]	25
19	Errno flags for <code>pton()</code> as defined in <code><errno.h></code> [8]	26
20	<code>ioctl()</code> flags defined in <code><bpf.h></code> [2]	42
21	Overview of <code>libnet_init()</code>	51
22	Overview of <code>libnet_build_udp()</code>	52
23	Overview of <code>libnet_build_ipv4()</code>	52
24	Overview of <code>libnet_build_ethernet()</code>	53
25	Overview of <code>libnet_autobuild_ipv4()</code>	54
26	Overview of <code>libnet_write()</code>	54
27	Overview of <code>libnet_name2addr4()</code>	55
28	Overview of <code>libnet_addr2name4()</code>	55
29	Overview of <code>libnet_toggle_checksum()</code>	56
30	Overview of <code>libnet_stats()</code>	56
31	Overview of <code>libnet_geterror()</code>	57
32	Overview of Network Layer Protocols	59
33	Overview of Transport Layer Protocols	59
34	Overview of <code>pcap_lookupdev()</code>	60
35	Overview of <code>pcap_lookupnet()</code>	61
36	Overview of <code>pcap_findalldevs()</code>	61
37	Overview of the <code>pcap_if_t</code> header	61
38	Overview of the <code>pcap_addr</code> header	61
39	Overview of the <code>sockaddr</code> header	62
40	Overview of <code>pcap_open_live()</code>	62
41	Overview of <code>pcap_compile()</code>	63
42	Overview of <code>pcap_setfilter()</code>	63
43	Overview of <code>pcap_next()</code>	64
44	Overview of <code>pcap_loop()</code>	64
45	Overview of callback	65
46	Common data link types	65
47	Protocol Types defined in <code><netinet/in.h></code> [8]	80
48	Linux Protocol Types defined in <code><linux/if_ether.h></code>	81
49	Socket level options for <code>setsockopt()</code> as defined in <code><errno.h></code> [8]	82
50	IP level options for <code>setsockopt()</code> as defined in <code><errno.h></code> [8]	83
51	Errno flags for <code>connect()</code> as defined in <code><errno.h></code> [8]	84
52	<code>ioctl()</code> flags defined in <code><bpf.h></code> [2]	85

List of Listings

1	sockaddr in cpp	26
2	sockaddr in6 cpp	26
3	sockaddr init cpp	26
4	sockaddr ll cpp	27
5	ifreq cpp	27
6	ifreq init cpp	27
7	mac cpp	36
8	bpf cpp	40
9	bpf2 cpp	41
10	bpf structs cpp	41
11	bpf structs filter cpp	42
12	rfc1071 checksum cpp	69
13	win socket cpp	70
14	libnet tcp c	72
15	sendpack c	75
16	dump c	76

1 Introduction

First of all we want to define the basic wordings and concepts used. Then the reader is given a short overview about the possibilities and restrictions the available APIs impose on the user.

Internet sockets are the common way to perform network communication implemented in most operating systems. They are usually provided by a socket API and are based upon the same principles as reading and writing a file. A program can get a socket via a function provided by the operating system. This function then returns a socket descriptor, usually a simple integer, similar to the ones provided by most operating systems for Read- and Write-Operations on files. This socket descriptor then can be used to write or read data from the socket. The data written to the socket is encapsulated by the operating system by adding headers and trailers and then the complete packet is sent via the network interface over the network to the target host. The data that has been received from the socket is presented to the program without the headers and trailers, only the user data is presented to the user. The sockets that work this way are the DATAGRAM- or STREAM-sockets provided via the Berkeley socket API. The user is unaware of the communication between the lower layers. All network communication steps, like for example the connection establishment, are taken care of without knowledge of the user. The user is only responsible for creating the socket and then providing the data that he wants to send to the correct function.

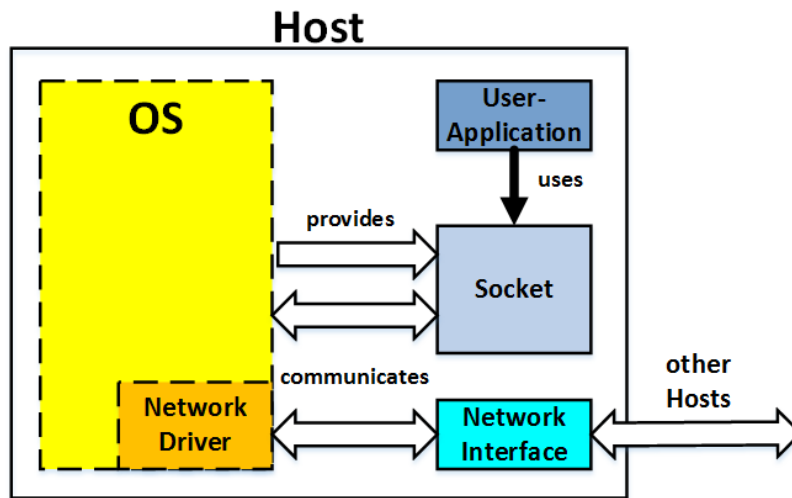


Figure 1: Socket provided by the operating system

The only parameters the user has to set are:

- The source socket address, a combination of IP address and port number can be set usually via the `bind()` function or defined directly with one of the `send()` functions.
- By choosing either the DATAGRAM- or STREAM-socket we can choose if we want to have a separate datagrams or a byte stream. This also chooses the Transport Layer Protocol that is being used (UDP or TCP).

If we want to gain access to the data of the lower layers we have several possibilities: RAW-sockets, PACKET-sockets, network drivers and Data Link layer APIs. The programming of Network Drivers will not be discussed further in this work, since we want to have a look at portable solutions that work for different operating systems. What can be accessed by the APIs we will present is shown in figure 2 on the following page.

With these APIs it is possible for an application to change and access the fields of the Network layers that are used for sending the data. This might be seen as a break with the traditional layering model, since we can influence the service the lower layers provide.

1.1 RAW-sockets

RAW-sockets are part of the standard Berkeley sockets and the socket API that is based upon it [1]. They are another option in addition to the already mentioned DATAGRAM- or STREAM-sockets to create data packets with the socket API [1]. In addition to simply sending data and defining address information the RAW-socket allows the user to access and manipulate the header and trailer information of the lower layers, more specifically with RAW-sockets to the Network and Transport layer (layer 3 and 4 of the OSI model). Since RAW-sockets are part of the Internet socket API they can only be used to generate and receive IP-Packets.

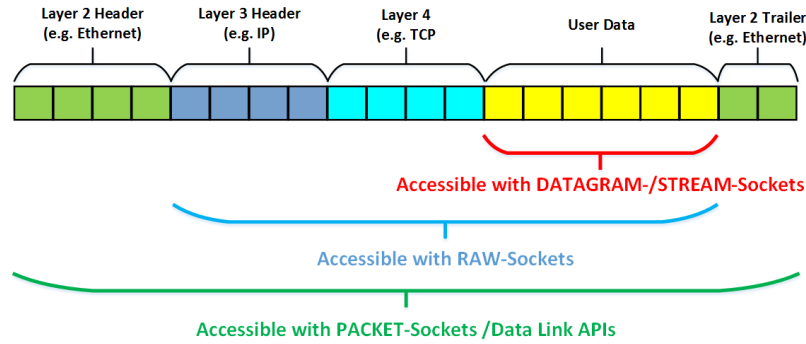


Figure 2: Overview over the layers and access possibilities

The biggest problem with RAW-sockets (also the PACKET-sockets we discuss later) is that there is no uniform API for using RAW-sockets under different operating systems:

- The provided APIs differ in regard to the used byte order. Depending on the operating system the fields of the packets have to be filled with data in network- or host-byte-order.
- Differences also exist in the types of packets and protocols that can be created using the RAW-socket API.
- The usage and functionality of the APIs is also different for each operating system.
- Some operating systems do not allow certain packet types to be received using the RAW-socket API.
- There are also differences in the definitions and paths of the necessary header files provided by each operating system.
- The required access levels for using the RAW-socket API can also differ. However most operating systems require root or admin access permissions to use them.

The user has to keep these things in mind if he wants to use the RAW-socket API, especially if he plans to use them across different operating systems.

1.2 PACKET-sockets and Data Link Layer APIs

If the user also wants to access the header and trailer of the Data Link layer then a Data Link Layer API has to be used. Under Linux (and Unix as well) this functionality is provided by the operating system as the so called PACKET-socket and can be seen as a special form of RAW-sockets, so the limitations and problems of the previous section also apply to them [1]. But since this is not a standardized functionality there might be a different API provided by the operating system to perform these operations. For example under BSD the so called BSD Packet Filter (BPF) not the PACKET-socket API provides access to the Data Link Layer [2]. In general all of these APIs usually allow access to the whole packet from the Data Link Layer (OSI Layer 2), the Network Layer (OSI Layer 3) up to the Transport layer (OSI Layer 4). Not all operating systems provide a simple interface to access the Data Link Layer. As an alternative we introduce a library that can be used to provide this functionality.

2 Implementations for different Operating Systems

First we want to give a short overview over some of the APIs and libraries that are available for RAW-socket and Data Link Layer network programming. In general we want to split the available libraries into two parts, the APIs provided by the operating systems and independent libraries that are available for multiple operating systems.

2.1 Windows

Windows as a operating system is pretty restricted when it comes to RAW-socket, PACKET-socket and Data Link Layer programming. In general it is difficult to get the programs running under Windows and the available options are pretty limited. For that reason in general 3rd party libraries are recommended if the user wants to do network programming on Windows systems and still write portable code.

2.1.1 Winsock-API

Winsock, currently available in version 2, is the Windows implementation of the Berkeley socket. Winsock is based on the implementation of sockets on Linux and BSD and allows both sending and receiving of RAW-Packets, but adds additional functionality.

Nevertheless, in the latest operating system versions (Windows XP and newer), there are the some restrictions the programmer has to consider according to the MSDN-library [3] for developers:

- Like in most operating system the Winsock-API also need admin access permits to work properly for RAW-sockets.
- A call to the bind function with a RAW-socket for the TCP protocol is not allowed, in fact TCP data cannot be sent over a RAW-socket at all.
- UDP datagrams with an invalid source address cannot be sent over RAW-sockets.
- The IP source address for any outgoing UDP datagram must exist on a network interface or the datagram is dropped. This change was made to limit the ability of malicious code to create distributed denial-of-service attacks and to limit the ability to send spoofed packets (TCP/IP packets with a forged source IP address).
- Especially receiving is much slower with RAW-sockets than using a customized driver based on the lower network layers, since all traffic is received (this consideration is the same for any other operating system).

2.1.2 winpcap

The ability of Winsock to also manipulate the Data Link Layer was removed in the current versions. Instead Windows does require to program a driver with the required functionality. With Network Driver Interface Specification (NDIS) version 6.40 (currently) the Operating system provides an API to do this if needed. This is now the preferred way to write own applications that want to access the Data Link Layer on Windows. Since want to only look at APIs and libraries provided by the operating system, we want to point to winpcap and libnet at this point as an alternative to the Winsock-API. Winpcap is the windows port of libpcap and allows generation and receiving of packets from the Data Link Layer upwards. Libnet is a network programming library available for many operating systems that allow easy sending of RAW-packets. All further details of winpcap and libnet can be found under their sections, if there are differences using the windows versions of the libraries they will be pointed out accordingly [4].

2.2 Linux

Linux is one of the operating systems for which it is easy to do RAW-socket and Data Link Layer Programming. It provides the APIs to do both, but the kernel has to be compiled with the option to support the options.

2.2.1 RAW-sockets

The RAW-socket is included in the socket-API as we already discussed. It allows both sending and receiving of RAW-packets, however the following items still have to be observed:

- Due to essential nature of header information for networking functionality and security using RAW-sockets requires root access permissions [5].
- The ports of the network layer are not endpoints anymore since RAW-sockets work on the layers below. Filtering based on ports has to be done manually [5].

- The `bind()` and `connect()` functions are no longer necessary [5]. `bind()` and `connect()` can still be used to define source address and target address to be entered automatically by the kernel [5]. Additionally a raw socket can be bound to a network device using `SO_BINDTODEVICE`.
- The `listen()` and `accept()` functions are without function, since the client-server semantic is no longer present [5]. When we use RAW-sockets we are sending unconnected packets [5].
- **IP-headers of RAW sockets can be manually created by the programmer if the option `IP_HDRINCL` is enabled [1]. This way, Raw sockets allow a programmer to implement new IP based protocols. If `IP_HDRINCL` is not enabled the IP header will be generated automatically.**
- When a RAW socket is created any IP based protocol can be specified [1]. This results in a socket only receiving messages of the type of the specified protocol.
- If a programmer does not want to specify a protocol when creating a RAW-socket he can also use the `IPPROTO_RAW` protocol (which implies that the headers will be created manually) [1]. This way he can send any IP based protocol. However this socket is not able to receive any IP packets, to receive all IP based packets a `PACKET`-socket has to be used.

2.2.2 PACKET-sockets

`PACKET`-sockets are a special type of RAW-sockets that allow access to the fields of the Data Link Layer. To use them there are some definitions required that not part of the socket-API, but are provided by the Linux Kernel. So while they are integrated in the socket-API of Linux, they are dependent on our operating system [5].

Similar restrictions than for the RAW-sockets have to be considered:

- Similar to RAW-sockets, due to essential nature of header information for networking functionality and security reasons using `PACKET`-sockets also requires root access permission [5].
- The ports of the Network layer are not endpoints anymore, since RAW-sockets work on the layers below [5]. Filtering based on ports has to be done manually.
- The `bind()` and `connect()` functions are no longer necessary [5]. `bind()` can still be used to define the interface over which the Data Link Layer packet should be sent [5]. `connect()` has no function anymore [5].
- The `listen()` and `accept()` functions are no use at all since the Client-Server-Semantic is no longer present [5]. When we use RAW-sockets we are sending unconnected packets.

2.3 Unix (FreeBSD, Mac OS X)

Like Linux under Unix operating systems there are also RAW-sockets provided with the Unix kernel. RAW-sockets are however more restricted than under Linux. Access to the Data Link Layer is possible in Unix via the Berkeley Packet Filters (BPF) provided by the operating system.

2.3.1 RAW-sockets

RAW-sockets are included in the socket-API in Unix like they are in Linux. In Unix there are different header files available and needed than in Linux, that makes some considerations necessary to keep code portable between these two operating systems [1]. Also there are some differences to the functionality.

Similar to Linux there are some restrictions to take into consideration:

- It is not possible to read packets for anything that has a handler (like TCP or UDP), but it is possible to read packets for other protocols like ICMP [6].
- For BSD and its ports it could be that the Packets are modified by the operating system before sending. For example with the release 10 the IP length is modified to the actual size of the IP header regardless of what is set by the programmer.
- Due to essential nature of header information for networking functionality and security reasons using RAW-sockets requires root access permission [5].
- The ports of the Network layer are not endpoints anymore, since RAW-sockets work on the layers below [5]. Filtering based on ports has to be done manually.
- The `bind()` and `connect()` functions are no longer necessary [5]. `bind()` and `connect()` can still be used to define source address and target address to be entered automatically by the kernel [5].

- The `listen()` and `accept()` functions are no use at all since the Client-Server-Semantic is no longer present [5]. When we use RAW-sockets we are sending unconnected packets.

2.3.2 Berkeley Packet Filter (BPF)

The Berkeley Packet Filter is the API provided by Unix operating systems to allow the user access to the Data Link Layer. It is compiled into the kernel of Unix-like hosts and allows access to all packets received at the Network Interface Controller (NIC) [7]. It has a built in filtering mechanism that allows the user to filter the received traffic for the packets he is interested in [7]. Since the RAW-socket implementation does not support receiving of packets that have a handler (like TCP or UDP) the Berkeley Packet Filter has to be used if the user wants to receive these packets [7].

2.4 OS independent

Other than the OS dependent APIs and the standard APIs provided by the standard Berkeley socket API there are also universal packet capturing and creation libraries. Of those we want to introduce two libraries, namely **libpcap** and **libnet** which are available for many operating systems (e.g. Windows, OS X, BSD, Linux) and therefore can make programming with RAW-sockets more portable. The libraries complement each other **pcap** is usually used to receive the packets while **libnet** provides easy mechanisms to send packets.

2.4.1 pcap

Pcap is an open library for packet capture. It allows the user to receive and filter all packets received at the network interface from the Data Link Layer upwards. It is available across different platforms, for Unix and Linux systems there exists libpcap (library for packet capture), the library for Windows is called Winpcap. It also supports injection of layer 2 packets, but the use of this functionality is very limited and not very comfortable. The use of the library might require admin or root permissions to work, depending on the operating system. One of the most well known programs that uses pcap is Wireshark.

2.4.2 libnet

Libnet is a library for constructing and sending packets from the Data Link Layer upwards. It is intended as the counterpart for pcap. What it does provide is a simple, modular interface for packet construction and injection. It gives programmers many helpful functions to make the construction of own packets very easy. Currently over 30+ protocols are supported by the library and it is available for many operating systems.

3 Programming with the APIs

3.1 General

We start by discussing some general topics that have to be considered for the development with any of the following APIs. Other than the topics presented here, it should be noted that the programmer also has to observe the particularities of the network protocols he wants to use. If these are not observed, like for example the correct order in that protocol information has to be exchanged, it cannot be ensured that we are able to communicate with the other host since the results than can be very different based on the implementation of the the target host.

3.1.1 Byte Order

In network communication the byte-order, also known as endianness, is essential for setting and interpreting the fields correctly. **It comes into play whenever a value does need more than 1 byte of storage space. Whenever this is the case we have to translate between network-byte-order and host-byte-order. The network-byte-order and host-byte-order may be different depending on the used protocol and operating system.** For the IP-protocol the network-byte-order is big-endian, so we have to translate whenever the byte-order of the host is not big-endian [5]. So we have two cases:

- When we send network packets all data that needs more than 1 byte it has to be translated from host-byte-order to network-byte-order.
- When we receive data on a socket all data that needs more than 1 byte it has to be translated from network-byte-order to host-byte-order.

For the transformations used for the IP-protocol the Berkeley socket API provides a set of standard functions which are available in the header file `<arpa/inet.h>`:

Function	Description
<code>uint32_t htonl(uint32_t hostlong)</code>	host-to-network
<code>uint16_t htons(uint16_t hostshort)</code>	host-to-network
<code>uint32_t ntohl(uint32_t netlong)</code>	network-to-host
<code>uint16_t ntohs(uint16_t netshort)</code>	network-to-host

Table 1: Byte-Order Transformation Functions [8]

With **Linux** we can be sure that the data receive/sent over a socket always is in network-byte-order [5].

For **Unix** it can be the case that some fields of the header are in host-byte-order depending of the release (for example BSD) [5].

3.1.2 Checksum

Many protocols used for network communication use a checksum to be able to detect transmission errors. There are different ways to calculate a checksum. If we choose to generate the packet headers completely by our self we have to calculate and set the applicable header fields manually. There are some points to take into consideration if we do this:

- Which algorithm has to be used to compute the checksum.
- Which fields of the header, sometimes even of multiple headers, have to be included in the calculation of the checksum [5]. Of course the user data might have to be included in the calculation as well.
- If we want to calculate the checksum all header information already has to be set with the correct data.
- Do we have to calculate the checksum or can it be calculated automatically by the network driver/hardware (for example the Ethernet checksum), the operating system or a predefined function call [5].

As a possible code example for the calculation of a checksum here is the listing for computing the Internet-checksum which was provided with the RFC 1071 which can be found in the appendix.

3.1.3 Type-Casting

One of the basic functionalities of the C programming language is that it supports Type-Casting. The ability to cast binary data to a data structure `struct` in C is essential to make working with header data easier for the programmer. To perform type casting we simply create a pointer to the desired structure:

```
struct header* protocolHeader;
```

Then we do need a pointer to the memory space where the binary data of the packet is stored, in our examples mostly a `char`-Array with a fixed `SIZE`[5]:

```
char buffer[SIZE];
```

Finally we can cast the binary data to our `struct` and afterwards are able to access the header fields[5]:

```
protocolHeader = (struct header*) buffer;
```

When accessing the header-field we still have to take care of the byte order like mentioned before. For the majority of the well known protocols there are predefined structs available in the header files included in the libraries. Some of the available headers are shown in the table below, but it does only show a part of the available header structures. More headers may be available depending on the operating system and the library.

Linux Under Linux typically the following header files are available:

Header-File	Description
<netinet/ip.h>	Defines macros, variables, and structures for IP.
<netinet/ip_icmp.h>	Defines macros, variables, and structures for ICMP.
<netinet/udp.h>	Defines macros, variables, and structures for UDP.
<netinet/tcp.h>	Defines macros, variables, and structures for TCP.
<netns/idp.h>	Defines IPX packet headers.
<netns/sp.h>	Defines SPX packet header.
<ssl.h>	Defines SSL prototypes, macros, variables, and structures.

Table 2: C-Header-Files for Network Headers [8]

3.1.4 Header-Positions

Another issue that has to be taken into consideration is the position of the packet headers in the binary data. The following has to be taken into consideration:

- Like it was shown in the beginning the header of a lower layer always encapsulate the header of the higher layers and the data.
- The header of the lowest layer can always be found at the beginning of the binary data packet. It might be the case that not all fields are present, for example this is the case for the Ethernet-Protocol in the Data Link layer. For the Ethernet-Protocol the preamble, start of frame delimiter and the frame checksum are removed by the network driver since they do belong more to the Physical Layer (OSI Layer 1)[5].
- If we type cast binary data to structs we have to add the length of the last header to the data pointer to get the beginning of the next header [5]:

```
next_header = (struct header*) (buffer + sizeof(struct header));
```

- To find the correct start position might be a bit more difficult than a simple add operation as shown here if the last header has a variable length with optional fields (like for example the IP-header or the TCP-header). In that case we must first get the fixed part of the header, then access the length field and finally compute the length of the optional header. Only after doing this we are able to find the start of the next header. If we are interested in the content of the optional fields we also have to access the header field containing the type of optional fields and cast it to the correct type accordingly.

3.2 RAW-sockets

The first API discussed here is the RAW-socket. It is available on Linux and Unix Systems. As already mentioned the PACKET-sockets can be seen as a specialized form of RAW-sockets and just allow access a lower layer then the normal RAW-socket. RAW-sockets allow access to the Network (OSI layer 4) and Internet (OSI layer 3) layer of the network stack [5]. The use of RAW-sockets is limited to processes with an effective user ID of 0 or the CAP_NET_RAW capability since they require **root**-access permissions [1]. In the following we will start by giving a generic overview over the functions available and necessary to create and work with RAW-sockets. After that we show the structure of different protocol level write and read operations.

3.2.1 socket ()

Both the read and write to a RAW-socket require the socket to be created first. For the creation of a socket the same function as for normal sockets is used. It is available in the `<netinet/in.h>` header and has the following form [8]:

```
int socket(int family, int type, int protocol)
```

The following parameters have to be defined:

- **family** expects a constant value that describes the used address family [8]. The following values are defined in `<sys/socket.h>`:

Constant	Description
AF_LOCAL	Local communication
AF_UNIX	Unix domain sockets
AF_INET	IP version 4
AF_INET6	IP version 6
AF_IPX	Novell IPX
AF_NETLINK	Kernel user interface device
AF_X25	Reserved for X.25 project
AF_AX25	Amateur Radio AX.25
AF_APPLETALK	Appletalk DDP
AF_PACKET	Low level packet interface
AF_ALG	Interface to kernel crypto API

Table 3: Address family constants provided in `<sys/socket.h>` [8]

The function only passes errors originating from the network to the user when the socket is connected. In that case only EMSGSIZE and EPROTO are passed for compatibility. If the IP_RECVERR flag is enabled, all network errors are saved in the error queue. It returns a non-negative socket-descriptor if it was created successful and -1 if an error occurred during creation of the socket. Additionally the variable `errno` defined in `<errno.h>` is set and can have the following values [8]:

Flag	Description
EACCES	User tried to send to a broadcast address without having the broadcast flag set on the socket.
EFAULT	An invalid memory address was supplied.
EINVAL	Invalid argument.
EMSGSIZE	Packet too big. Either Path MTU Discovery is enabled or the packet size exceeds the maximum allowed IPv4 packet size of 64KB.
EOPNOTSUPP	Invalid flag has been passed to a socket call.
EPERM	The user doesn't have permission to open raw sockets.
EPROTO	An ICMP error has arrived reporting a parameter problem.

Table 4: Errno flags for `socket ()` as defined in `<errno.h>` [8]

Users have to be aware that setting the **family** of course also influences which **protocol** can be chosen later. The usual option to use with RAW-sockets is the `AF_INET` to send and receive IP version 4 packets.

- **type** defines the socket type. The following values are defined in `<sys/socket.h>`:

Constant	Description
SOCK_STREAM	Stream (connection) socket
SOCK_DGRAM	Datagram (connection-less) socket
SOCK_RAW	RAW socket
SOCK_RDM	Reliably-delivered message
SOCK_SEQPACKET	Sequential packet socket
SOCK_PACKET	Linux specific way of getting packets at the dev level.

Table 5: The socket types are defined in `<sys/socket.h>`[8]

Since we want to work with RAW-sockets we only use the `SOCK_RAW` constant in the following sections. From Linux kernel 2.6.27, there is a second purpose for the type argument, it additionally allows to modify the behavior of the socket by including the following options with bit-wise OR [8]:

- **SOCK_NONBLOCK**: Sets the `O_NONBLOCK` file status flag on the new open socket descriptor. Using this changes the socket to a non-blocking one.
- **SOCK_CLOEXEC**: Sets the close-on-exec flag for the new socket descriptor. Useful if the program creates child processes which use `exec()` to run another program. In that case this flag will prevent the other program from using this socket descriptor.

- **protocol** defines like the name implies the protocol of the packets that can be sent and received with the socket [8]. The protocol numbers are defined by the IANA (Internet Assigned Numbers Authority) and a complete list can be found on their website. The table ?? shows the constants for protocols which are defined in the `<netinet/in.h>` header file.

Again, users have to be aware of which **family** was chosen for the first option, since only protocols of that family can be chosen as **protocol**. So with our usual `AF_INET` option selected we can only use IP-based protocols [8].

Also it should already be noted here that the `IPPROTO_RAW` constant in for most operating system (depends on the Linux/Unix distribution) does imply that the socket also expects an IP header to be manually created by the user [8]. If we chose this constant we have layer 3 write access as a result. The regular way to do this is to use another constant and then to change the socket options with `setsockopt()` as described in section 3.2.3 on the next page.

In addition to these constants we could also use constants that are defined for layer 2 (PACKET-sockets) to access their protocol information. These constants are operating system dependent and therefore code that uses them cannot be ported as easy as the general constants from the previous table. As an example, table 48 in the appendix shows the constants in Linux for Ethernet protocols which are defined in the `<linux/if_ether.h>` header.

Unix The `socket()` function is POSIX and 4.4BSD conform [8]. The `SOCK_NONBLOCK` and `SOCK_CLOEXEC` flags are Linux-specific [8]. The function is generally portable to BSD systems, but some systems require the `<sys/types.h>` header to work [8]. The manifest constants might also be different under some BSD versions.

3.2.2 setsockopt()

The `setsockopt()` function like the name implies can be used to change the options that are selected for the socket. The function can manipulate the options for different protocol levels such as IP or TCP, but also for the sockets level API by setting the level to `SOL_socket`. The function which is defined in the `<sys/socket.h>` header has the following form [8]:

```
int setsockopt(int sockfd, int level, int optname, const void * optval, socklen_t
optlen)
```

The following parameter can be set [8]:

- **sockfd** - Specifies the socket for which the options should be set.
- **level** - The protocol level of the option we want to set.

- **optname** - The name of the option we want to set. Together with the **optval** and **optlen** it is passed uninterpreted to the protocol module for processing.
- **optval** - The buffer for the option we want to specify, usually an Integer. It should then be non-zero to enable a Boolean option and zero to disable it.
- **optlen** - The length of the buffer that **optval** points to in bytes.

We want to only give an short overview over the most interesting options for us, the socket level and the IP protocol level, which can be set with the `setsockopt()` function. Additional information about the use and meaning of the options can be found in the programmer manuals for the protocols. They exist for all protocols available for the respective operating system or distribution.

For the socket level API that can be accessed with `SOL_socket` as the level, the options in table 49 in the appendix are supported and already included in the `<sys/socket.h>` header.

For the IP level API that can be accessed with `IPPROTO_IP` as the level, table 50 in this appendix shows the following options supported and already included in the `<netinet/ip.h>` header.

The function returns zero on success and -1 on an error [8]. Additionally the variable `errno` defined in `<errno.h>` is set and can have the following values:

Flag	Description
EBADF	The argument <code>sockfd</code> is not a valid descriptor.
EFAULT	The address pointed to by <code>optval</code> is not in a valid part of the process address space.
EINVAL	<code>optlen</code> invalid or invalid <code>optval</code> .
ENOPROTOOPT	The option is unknown at the level indicated.
ENOTSOCK	The argument <code>sockfd</code> is a file, not a socket.

Table 6: Errno flags for `setsockopt()` as defined in `<errno.h>`
[8]

Unix The function is POSIX and 4.4BSD conform [8]. According to the POSIX-standard the use of `setsockopt()` does only require to include `<sys/socket.h>` header, so under Linux it can be used by only including this single header [8]. To use it under Unix also require the `<sys/types.h>` header to be included [8]. For portability it is generally recommended to include both.

3.2.3 `getsockopt()`

For retrieving an specified option of the socket the `getsockopt()` function defined in the `<sys/socket.h>` header can be used [8]:

```
int getsockopt(int sockfd, int level, int optname, void * optval, socklen_t *
optlen)
```

The following parameter have to be set [8]:

- **sockfd** - Specifies the socket for which the options should be retrieved.
- **level** - The protocol level of the option we want to retrieve.
- **optname** - The name of the option we want to retrieve. Together with the **optval** and **optlen** it is passed uninterpreted to the protocol module for processing.
- **optval** - The buffer for the option we want to retrieve.
- **optlen** - The length of the buffer that **optval** points to in bytes. Initially it should contain the length of the buffer, after the call it indicates the actual size of the value returned. In case no value is supplied or returned it might return NULL.

The function returns zero on success and -1 on an error [8]. Additionally the variable `errno` defined in `<errno.h>` is set and can have the following values:

Flag	Description
EBADF	The argument <code>sockfd</code> is not a valid descriptor.
EFAULT	The address pointed to by <code>optval</code> or <code>optlen</code> is not in a valid part of the process address space.
EINVAL	<code>optlen</code> invalid or invalid <code>optval</code> .
ENOPROTOOPT	The option is unknown at the level indicated.
ENOTSOCK	The argument <code>sockfd</code> is a file, not a socket.

Table 7: Errno flags for `getsockopt()` as defined in `<errno.h>` [8]

Unix The function is POSIX and 4.4BSD conform [8]. According to the POSIX-standard the use of `getsockopt()` does only require to include `<sys/socket.h>` header, so under Linux it can be used by only including this single header [8]. To use it under Unix also require the `<sys/types.h>` header to be included [8]. For portability it is generally recommended to include both.

3.2.4 `bind()`

After creating a socket like discussed in the previous section 3.2.1 on page 15, we can bind the created socket to a specific address. Traditionally this is also called *assigning a name to a socket*. For RAW-sockets and PACKET-sockets this is optional, but we use it to define the source address of our packets with it and also define from which network-interface we want to read packets. Other socket types might require to be bound to a specific address before they can be used. For binding the socket to an IP-address we can use the `bind()` function which is defined in the `<sys/socket.h>` header and has the following form [8]:

```
int bind(int sockfd, const struct sockaddr * addr, socklen_t addrlen)
```

The following parameters can be set [8]:

- **sockfd** - Specifies the socket for which the address should be set.
- **addr** - The address structure containing the address to be set.
- **addrlen** - The length of the address structure in bytes.

The function returns zero on success and -1 on an error [8]. Additionally the variable `errno` defined in `<errno.h>` is set and can have the following values:

Flag	Description
EACCES	The address is protected, and the user is not the superuser.
EADDRINUSE	The given address is already in use.
EADDRINUSE	All port numbers in the port range are currently in use.
EBADF	<code>sockfd</code> is not a valid file descriptor.
EINVAL	The socket is already bound to an address.
EINVAL	<code>addrlen</code> is wrong, or <code>addr</code> is not a valid address for this socket's domain.
ENOTSOCK	The file descriptor <code>sockfd</code> does not refer to a socket.

Table 8: Errno flags for `bind()` as defined in `<errno.h>` [8]

Unix According to the POSIX-standard the use of `bind()` does only require to include `<sys/socket.h>` header, so under Linux it can be used by only including this single header [8]. To use it under Unix also requires the `<sys/types.h>` header to be included [8]. For portability it is generally recommended to include both. For a UNIX domain (AF_UNIX) the following `errno` are additionally defined [8]:

Flag	Description
EADDRNOTAVAILA	Nonexistent interface was requested or the requested address was not local.
ENAMETOOLONG	<code>addr</code> is too long.
ENOMEM	Insufficient kernel memory was available.

Table 9: Errno flags for UNIX-`bind()` as defined in `<errno.h>` [8]

3.2.5 getsockname()

To get the currently defined name of a socket the `getsockname()` function can be used. It is defined in the `<sys/socket.h>` and has the following form [8]:

```
int getsockname(int sockfd, struct sockaddr * addr, socklen_t * addrlen)
```

The following parameters can be set [8]:

- **sockfd** - Specifies the socket for which the address should be retrieved.
- **addr** - The address structure for the returned address that is set for the socket. The returned address is truncated if the buffer is too small.
- **addrlen** - The length of the address structure in bytes. It should be initialized to the size of the buffer pointed to by **addr**. Contains the actual size of the socket address after return.

The function returns zero on success and -1 on an error [8]. Additionally the variable `errno` defined in `<errno.h>` is set and can have the following values:

Flag	Description
EBADF	The argument <code>sockfd</code> is not a valid file descriptor.
EFAULT	The <code>addr</code> argument points to memory not in a valid part of the process address space.
EINVAL	<code>addrlen</code> is invalid.
ENOBUFS	Insufficient resources were available in the system to perform the operation.
ENOTSOCK	The file descriptor <code>sockfd</code> does not refer to a socket.

Table 10: Errno flags for `getsockname()` as defined in `<errno.h>` [8]

Unix The function is compliant to POSIX and 4.4BSD [8]. In some distributions the third argument is in reality a pointer to an `int`[8]. This could still be the case in some distributions.

3.2.6 connect()

This function can be used to initiate a connection to a specific destination host and is required for the `write()`, `send()`, `read()` and `recv()` functions. For connection based protocols like `SOCK_STREAM` this function also will try to connect to the host on the other side [8]. For Datagram based protocols only the default destination is defined with this function [8]. To use this function the header `<sys/socket.h>` (for the function) has to be included. The function is called like this [8]:

```
int connect(int sockfd, const struct sockaddr * addr, socklen_t addrlen)
```

With the following parameters that can be set [8]:

- **sockfd** - Specifies the socket for which the address should be set.
- **addr** - The function specifies a `struct sockaddr *` with this parameter that contains the socket address family and the protocol address for that family. With this the address to connect to can be defined.
- **addrlen** - Defines how long the address that is passed to the function is.

The function returns zero on success and -1 on an error [8]. Additionally the variable `errno` defined in `<errno.h>` is set and can have the values shown in table 51 in the appendix.

Unix The `connect()` function might require the `<sys/types.h>` header for some BSD systems [8]. Also the third argument might be an `int` depending on the system version [8].

3.2.7 Read

We can access the socket with three different function to retrieve data from it. For the `read()` and `recv()` usually the source address is defined by connecting the socket to a specific host. This can be done by using the `connect()` function as described in section 3.2.6.

read() The read function works identical to the `read()` on a file. As already said the the socket can be connected to a specific host first by using the `connect` function as described in section 3.2.6 on the preceding page. For the `read()` function the definitions can be found in the `<unistd.h>` header. The function has the following form [8]:

```
int read(int fd, char * Buff, int NumBytes)
```

The following parameters have to be set [8]:

- **fd** - Takes a file descriptor from which the data should be read, for RAW-sockets we can use the socket descriptor instead.
- **Buff** - Defines the memory space for the binary data we want to read from the socket.
- **NumBytes** - How many bytes should be read from the socket, usually initialized with the size of our memory space.

A usual call of the read function for our purposes would look like this:

```
ssize_t = read(socket, packet, sizeof(packet));
```

On success the number of bytes is returned and on an error -1 [8]. Additionally the variable `errno` defined in `<errno.h>` is set and can have the following values:

Flag	Description
EAGAIN	The file descriptor <code>fd</code> refers to a file other than a socket and has been marked non-blocking (<code>O_NONBLOCK</code>) and the read would block.
EWOULDBLOCK	The file descriptor <code>fd</code> refers to a socket and has been marked non-blocking (<code>O_NONBLOCK</code>) and the read would block.
EBADF	<code>fd</code> is not a valid file descriptor or is not open for reading.
EFAULT	<code>buf</code> is outside your accessible address space.
EINTR	The call was interrupted by a signal before any data was read.
EINVAL	<code>fd</code> is attached to an object which is unsuitable for reading or the file was opened with the <code>O_DIRECT</code> flag.
EIO	I/O error.

Table 11: Errno flags for `read()` as defined in `<errno.h>` [8]

The possible return values and error flags can be found in section 13 on the next page.

recv() The `recv()` function is another possibility to retrieve data from a socket and does not require an address to be defined as well. As already said the the socket can be connected to a specific host first by using the `connect()` function as described in section 3.2.6 on the preceding page. To use the `recv()` function the headers `<sys/types.h>` (for the data types) and `<sys/socket.h>` (for the function) have to be included. The function is called like this [8]:

```
ssize_t recv(int sockfd, void * buf, size_t len, int flags)
```

The `recv()` function normally is used on a connected socket (after connecting with the `connect()` function), but also works with RAW-sockets. Similar parameters to the `read` function have to be set [8]:

- **sockfd** - Specifies the socket from which data should be read.
- **buf** - Defines the memory space for the binary data we want to read from the socket.
- **len** - How many bytes should be read from the socket, usually initialized with the size of our memory space.
- **flags** - The flags that can be set for the function are listed in table 12 on the next page.

recvfrom() The `recvfrom()` function also allows us to define the address of a host we want to receive data from. To use this function the headers `<sys/types.h>` (for the data types) and `<sys/socket.h>` (for the function) have to be included. The function is called like this [8]:

```
ssize_t recvfrom(int sockfd, void * buf, size_t len, int flags, struct sockaddr *
src_addr, socklen_t * addrlen)
```

The parameters are identical to the `recv()` function, it only has two additional parameters [8]:

- **sockfd** - Specifies the socket from which data should be read.
- **buf** - Defines the memory space for the binary data we want to read from the socket.
- **len** - How many bytes should be read from the socket, usually initialized with the size of our memory space.
- **src_addr** - The function returns a `struct sockaddr *` with this parameter that contains the socket address family and the protocol address for that family. The address is filled in by the called protocol and contains the source address of the packet received. The address will be truncated in case it is too long for the provided buffer. The parameter can be set to `NULL`, then it will not be filled by the protocol.
- **addrlen** - The length of the address. In case the **src_addr** parameter was too small for the returned address an address length greater than the one provided to the function will be returned.
- **flags** - The flags that can be set for the function are listed in table 12.

The possible return values and error flags can be found in section 13.

flags and errors for the `recv()`- and `recvfrom()`-functions Here are the possible flags for the `recv()` and `recvfrom()` functions:

Flag	Description
MSG_DONTWAIT	Enables non-blocking operation, if the operation would block the call fails with the error <code>EAGAIN</code> or <code>EWOULDBLOCK</code> .
MSG_ERRQUEUE	Specifies that queued errors should be received from the socket error queue, the buffer has to be supplied by the user.
MSG_OOB	This flag requests receipt of out-of-band data that would not be received in the normal data stream.
MSG_PEEK	This flag causes the receive operation to return data from the beginning of the receive queue without removing that data from the queue.
MSG_TRUNC	For raw (<code>AF_PACKET</code>), Internet datagram, netlink and UNIX datagram: return the real length of the packet or datagram even when it was longer than the passed buffer.
MSG_WAITALL	This flag requests that the operation block until the full request is satisfied. However the call may still return less data than requested if a signal is caught, an error or disconnect occurs or the next data to be received is of a different type than that returned.

Table 12: Flags for `recv()` and `recvfrom()` defined in `<sys/socket.h>` [8]

The functions `recv` and `recvfrom()` returns the count of bytes read on success and -1 on an error [8]. Additionally the variable `errno` defined in `<errno.h>` is set and can have the values shown in the following table:

Flag	Description
<code>EAGAIN</code>	The socket is marked non-blocking and the receive operation would block, or a receive timeout had been set and the timeout expired before data was received.
<code>EWOULDBLOCK</code>	The socket is marked non-blocking and the receive operation would block, or a receive timeout had been set and the timeout expired before data was received.
<code>EBADF</code>	The argument <code>sockfd</code> is an invalid descriptor.
<code>EFAULT</code>	The receive buffer pointer(s) point outside the processes address space.
<code>EINTR</code>	The receive was interrupted by delivery of a signal before any data were available.
<code>EINVAL</code>	Invalid argument passed.
<code>ENOTSOCK</code>	The argument <code>sockfd</code> does not refer to a socket.

Table 13: `Errno` flags for `recv()` and `recvfrom()` as defined in `<errno.h>` [8]

Unix The `read()` function is POSIX and 4.3BSD conform and should be available as described [8]. The `recv()` and `recvfrom()` are also conform to POSIX, but only the `MSG_OOB`, `MSG_PEEK` and `MSG_WAITALL` flags are described in the standard [8]. It is also conform to 4.4BSD, but there might be differences in the data types depending on the library used in the system [8].

3.2.8 Write

We can send data over the socket with three different functions. The `write()` and `send()` functions we discuss in the beginning both require a destination address to be defined first. This can be done by using the `connect()` function described section 3.2.6 on page 19. The last function we discuss, `sendto()`, has the option to define the destination address, but it can also be set by using the `connect()` function. For all functions we must take the buffer constraints into consideration, since otherwise the write connection will be closed before all data is transferred completely.

write() Identical to how we use a the `read()` function to get data from a socket descriptor instead of a file descriptor, we can also use the `write()` function that is usually used for data output to a file for sending data over a socket. As already mentioned to use `write()` with a socket a valid destination address has to be provided to the socket first for it to work. That can be done by using the `connect()` function which is discussed in section 3.2.6 on page 19. The `write()` function is defined in `<unistd.h>` and has the following from [8]:

```
ssize_t write(int fd, const void * buf, size_t count)
```

The parameters are:

- **fd** - In our case specifies the socket to which we want to write data.
- **buf** - Defines the memory space for the binary data we want to write.
- **count** - How many bytes should be written from the buffer to the socket. If there is not enough free space on the physical medium used, less bytes might be written. Also if the process is interrupted by a signal, there also might be less bytes written than specified.

The function returns the amount of bytes written on success and -1 on an error [8]. Additionally the variable `errno` defined in `<errno.h>` is set and can have the following values:

Flag	Description
EAGAIN	The file descriptor refers to a file other than a socket and has been marked non-blocking, and the write would block.
EWOULDBLOCK	The file descriptor <code>fd</code> refers to a socket and has been marked non-blocking (<code>O_NONBLOCK</code>), and the write would block.
EBADF	<code>fd</code> is not a valid file descriptor or is not open for writing.
EDESTADDRREQ	<code>fd</code> refers to a datagram socket for which a peer address has not been set using <code>connect()</code> .
EFAULT	<code>buf</code> is outside your accessible address space.
EINTR	The call was interrupted by a signal before any data was written.
EINVAL	<code>fd</code> is attached to an object which is unsuitable for writing or the file was opened with the <code>O_DIRECT</code> flag and either the address specified in <code>buf</code> , the value specified in <code>count</code> or the current file offset is not suitably aligned.
EIO	A low-level I/O error occurred while modifying the inode.
ENOSPC	The device containing the file referred to by <code>fd</code> has no room for the data.
EPIPE	<code>fd</code> is connected to a pipe or socket whose reading end is closed.

Table 14: Errno flags for `write()` as defined in `<errno.h>` [8]

send() The `send()` function is another function to send data over a socket. Like `write()` it does require that a destination address is specified first before it can be used, since it does require the socket to be in a connected state [8]. This can be achieved by using the `connect()` function discussed in section 3.2.6 on page 19. There is no indication of a failure implicitly shown when using `send()`, only locally detected errors are indicated by returning -1 [8]. If `send()` is used to transmit messages, then it will block (if not set otherwise) if a message does not fit in the buffer [8]. To use this function the headers `<sys/types.h>` (for the data types) and `<sys/socket.h>` (for

the function) have to be included. The function is called like this [8]:

```
ssize_t send(int sockfd, void * buf, size_t len, int flags)
```

With the following parameters that can be set [8]:

- **sockfd** - Specifies the socket over which the data should be sent.
- **buf** - Defines the memory space for the binary data we want to send over the socket.
- **len** - How many bytes should be read from the memory space given by the parameter buffer.
- **flags** - The flags that can be set for the function are listed in table 15 on the next page

Like already mentioned, the function only indicates local errors by setting the return value to -1 on an error and returning the number of characters (bytes) sent on success [8]. Additionally the variable `errno` defined in `<errno.h>` is set and can have the values described in the table 16 on the following page. Other errors might be reported by the used protocol modules.

sendto() In comparison to the possibilities to send data we discussed so far, the `sendto()` function allows us to define a specific address to which the data should be sent without having to call `connect()` first to set the destination address [8]. It should be noted that if the `sendto()` function is used on a connection-mode socket the additional address information is ignored and an `EISCONN` error might be returned in `errno`, if they are not set to `NULL` and zero respectively [8]. This is the case since for such a socket the destination is already implied by the connection type. To use this function the headers `<sys/types.h>` (for the data types) and `<sys/socket.h>` (for the function) have to be included. The function is called like this [8]:

```
ssize_t sendto(int sockfd, void * buf, size_t len, int flags, struct sockaddr *  
dest_addr, socklen_t addrlen)
```

With the following parameters that can be set [8]:

- **sockfd** - Specifies the socket over which the data should be sent.
- **buf** - Defines the memory space for the binary data we want to send over the socket.
- **len** - How many bytes should be read from the memory space given by the parameter buffer.
- **dest_addr** - The function defines a `struct sockaddr *` with this parameter that contains the socket address family and the protocol address for that family. The parameter can be set to `NULL`, then it will not be filled in.
- **addrlen** - The length of the address provided to the socket.
- **flags** - The flags that can be set for the function are listed in table 15 on the next page.

Like already mentioned, the function only indicates local errors by setting the return value to -1 on an error and returning the number of characters sent on success [8]. Additionally the variable `errno` defined in `<errno.h>` is set and can have the values described in the table 16 on the following page. Other errors might be reported by the used protocol modules.

Flags and errnos for the `send()` and `sendto()` functions The following flags can be set for both the `send()` and `sendto()` functions:

Flag	Description
MSG_CONFIRM	Only valid on SOCK_DGRAM and SOCK_RAW. Tell the layer 2 that you got a successful reply from the other side.
MSG_DONTROUTE	Don't use a gateway to send out the packet, only send to hosts on directly connected networks.
MSG_DONTWAIT	Enables non-blocking operation, if the operation would block EAGAIN or EWOULDBLOCK is returned.
MSG_EOR	Terminates a record (when this notion is supported).
MSG_MORE	The caller has more data to send. This flag is used with UDP/TCP sockets.
MSG_NOSIGNAL	Requests not to send SIGPIPE on errors on stream oriented sockets when the other end breaks the connection. The EPIPE error is still returned.
MSG_OOB	Sends out-of-band data on sockets that support this notion, the underlying protocol must also support out-of-band data.

Table 15: Flags for `send()` and `sendto()` as defined in `<errno.h>` [8]

The following `errno` constants are defined in `<errno.h>` for both the `send()` and `sendto()` functions:

Flag	Description
EAGAIN	The socket is marked non-blocking and the requested operation would block.
EWOULDBLOCK	The socket is marked non-blocking and the requested operation would block.
EBADF	An invalid descriptor was specified.
ECONNRESET	Connection reset by peer.
EDESTADDRREQ	The socket is not connection-mode, and no peer address is set.
EFAULT	An invalid user space address was specified for an argument.
EINTR	A signal occurred before any data was transmitted.
EINVAL	Invalid argument passed.
EISCONN	The connection-mode socket was connected already but a recipient was specified.
EMSGSIZE	The socket type requires that message be sent atomically, and the size of the message to be sent made this impossible.
ENOBUFS	The output queue for a network interface was full. This generally indicates that the interface has stopped sending, but may be caused by transient congestion.
ENOMEM	No memory available.
ENOTCONN	The socket is not connected, and no target has been given.
ENOTSOCK	The argument <code>sockfd</code> is not a socket.
EOPNOTSUPP	Some bit in the flags argument is inappropriate for the socket type.
EPIPE	The local end has been shut down on a connection oriented socket. In this case the process will also receive a SIGPIPE unless MSG_NOSIGNAL is set.

Table 16: Errno flags for `send()` and `sendto()` as defined in `<errno.h>` [8]

Unix The `write()` function is POSIX and 4.3BSD compliant, so it should work in BSD based systems without any changes [8]. The `send()` and `sendto()` are also conform to POSIX, but only the MSG_OOB, MSG_PEEK and MSG_WAITALL flags are described in the standard [8]. It is also conform to 4.4BSD, but there might be differences in the data types depending on the library used in the system [8].

3.2.9 `close()`

After we are finished with our network communication, we can close the socket like we would do it for a file. The `close()` function is defined in the `<unistd.h>` and looks like this [8]:

```
int close(int fd)
```

The only parameter is the socket descriptor that should be closed [8]. All locks held by the associated process are released as well [8]. The function returns zero on success and -1 on an error [8]. Additionally the variable `errno` defined in `<errno.h>` is set and can have the following values:

Flag	Description
EBADF	fd is not a valid file descriptor.
EINTR	The <code>close()</code> call was interrupted by a signal.
EIO	An I/O error occurred.

Table 17: Errno flags for `close()` as defined in `<errno.h>` [8]

Unix Since the `close()` function is part of the POSIX-standard there is no difference to the call between Linux-distributions and Unix-systems [8].

3.2.10 `inet_ntop()`

`inet_ntop()` is a helper function to convert from a network address to a human readable character string. It currently supports IP version 4 and version 6 addresses [8]. The function is defined in the `<arpa/inet.h>` and looks like this [8]:

```
const char * inet_ntop(int af, const void * src, char *dst, socklen_t size)
```

With the following parameters that can be set [8]:

- **af** - The address family we want to convert from. Currently only `AF_INET` for IP version 4 and `AF_INET6` for IP version 6 addresses are supported.
- **src** - The network address structure we want to convert.
- **dst** - A pointer to a buffer for the resulting character string.
- **size** - The size of the buffer.

On success the `dst` parameter will contain the human readable form of the address [8]. If an error occurs `NULL` is returned [8]. Additionally the variable `errno` defined in `<errno.h>` is set and can have the following values:

Flag	Description
EINVAL	The parameter af was not a valid address family.
ENOSPC	The converted address string would exceed the size given by <code>size</code> .

Table 18: Errno flags for `ntop()` as defined in `<errno.h>` [8]

Unix The function is a part of the POSIX standard and should work as described on all Linux and Unix systems [8].

3.2.11 `inet_pton()`

`inet_pton()` is a helper function to convert from a human readable character string to a network address. It currently supports IP version 4 and version 6 addresses [8]. The function is defined in the `<arpa/inet.h>` and looks like this [8]:

```
int inet_pton(int af, const char * src, void *dst)
```

With the following parameters that can be set [8]:

- **af** - The address family we want to convert from. Currently only `AF_INET` for IP version 4 and `AF_INET6` for IP version 6 addresses are supported.
- **src** - The character string we want to convert. For the `AF_INET` address family the preferred format is the dotted-decimal format, "ddd.ddd.ddd.ddd". For the `AF_INET6` address family the preferred format is the hexadecimal format "x:x:x:x:x:x:x:x", but the usual abbreviation forms are also supported.
- **dst** - A pointer to the resulting network address structure.

On success the `dst` parameter will contain the network address structure and the return value is set to 1 [8]. If the `src` parameter did not contain a valid address string 0 is returned [8]. If an error occurred -1 is returned and additionally the variable `errno` defined in `<errno.h>` is set and can have the following values:

Flag	Description
EAFNOSUPPORT	The parameter af was not a valid address family.

Table 19: Errno flags for `pton()` as defined in `<errno.h>` [8]

Unix The function is a part of the POSIX standard and should work as described on all Linux and Unix systems [8].

3.2.12 Data Types

There are several structures that we use for all system calls and functions we use when working with network functions. Here we want to show these structures and show a way to initialize them. The first structs we want to introduce are the `sockaddr` structs that are used to hold a layer 3 address and are defined in the `<netinet/in.h>` header. The struct `sockaddr` is the basic definition of an address that many of the functions discussed so far take as a parameter. It is possible to cast a pointer to one of these structures into each other without any harm [9].

See listing ?? in the appendix.

The struct `sockaddr_in` is the struct we use to hold an IP version 4 address. Note the `sin_zero` field, for which it is sometimes recommended that it should be set to zero, even though it is not even mentioned in the Linux programmer's manual [9]. We recommend to initialize the whole struct with 0 using the `memset()` function. The other fields are usual information for IP version 4 addresses.

Listing 1: `sockaddr` in cpp

```

1 struct sockaddr_in {
2     short      sin_family;   // e.g. AF_INET, AF_INET6
3     unsigned short sin_port; // e.g. htons(3490)
4     struct in_addr sin_addr; // see struct in_addr, below
5     char       sin_zero[8];  // zero this if you want to
6 };

```

The struct `sockaddr_in6` is the struct we use to hold an IP version 6 address. Similar to the struct for the IP version 4 address it contains all the required fields for the protocol information.

Listing 2: `sockaddr_in6` in cpp

```

1 struct sockaddr_in6 {
2     u_int16_t    sin6_family; // address family, AF_INET6
3     u_int16_t    sin6_port;   // port number, Network Byte Order
4     u_int32_t    sin6_flowinfo; // IPv6 flow information
5     struct in6_addr sin6_addr; // IPv6 address
6     u_int32_t    sin6_scope_id; // Scope ID
7 };

```

The following example adapted from [9] shows exemplary how to init an IP version 4 address:

Listing 3: `sockaddr` init cpp

```

1 // the struct for the socket-address
2 struct sockaddr_in ip4addr;
3
4 // init the address family to internet addresses
5 ip4addr.sin_family = AF_INET;
6
7 // fill out the port
8 ip4addr.sin_port = htons(3490);
9
10 // fill in the ethernet address
11 inet_pton(AF_INET, "10.0.0.1", &ip4addr.sin_addr);

```

There is also a structure we can use to define a layer 2 address, for example a MAC-address or an index for one of the network interfaces of the host. It is available in the `<linux/if_packet.h>` header.

Listing 4: sockaddr ll cpp

```
1 struct sockaddr_ll {
2     unsigned short sll_family;    // family is always AF_PACKET
3     unsigned short sll_protocol; // physical layer protocol
4     int sll_ifindex;             // interface number
5     unsigned short sll_hatype;   // header type
6     unsigned char sll_pkttype;   // packet type
7     unsigned char sll_halen;     // length of the address
8     unsigned char sll_addr[8];   // physical layer address
9 };
```

Many `ioctl()` calls we use to get network device information return an `ifreq` structure [10]. We usually specify which device to affect by setting `ifr_name` to the name of the interface [10]. It is defined in the `<net/if.h>` header.

Listing 5: ifreq cpp

```
1 struct ifreq {
2     char ifr_name[IFNAMSIZ]; /* Interface name */
3     union {
4         struct sockaddr ifr_addr;
5         struct sockaddr ifr_dstaddr;
6         struct sockaddr ifr_broadaddr;
7         struct sockaddr ifr_netmask;
8         struct sockaddr ifr_hwaddr;
9         short ifr_flags;
10        int ifr_ifindex;
11        int ifr_metric;
12        int ifr_mtu;
13        struct ifmap ifr_map;
14        char ifr_slave[IFNAMSIZ];
15        char ifr_newname[IFNAMSIZ];
16        char *ifr_data;
17    };
18 };
```

One example is the use of the `ifreq` for getting the device index of an network interface as shown in the following code:

Listing 6: ifreq init cpp

```
1 // struct for the interface definition
2 struct ifreq ifr;
3
4 // the name of the device we want to use
5 char * device = "wlan0";
6
7 // set the structs to zero
8 bzero(&ifr , sizeof(ifr));
9
10 // define the name of the interface we want to set
11 strncpy((char *)ifr.ifr_name ,device , IFNAMSIZ);
12
13 // get device index
14 ioctl(sockfd , SIOCGIFINDEX , &ifr)
```

3.2.13 Layer 4

Now we can start with describing the different layers we can access and write with the RAW-sockets. We start by discussing working with RAW-sockets for Transport layer (OSI layer 4) access. The basic operation is generally very similar for all layers we want to work with. For example the basic call to create a RAW-socket for IP looks like this [5]:

```
sockfd = socket(AF_INET, SOCK_RAW, int protocol )
```

In this example the `AF_INET` specifies IP as the protocol, the `SOCK_RAW` defines the socket type as a RAW-Socket and the `protocol` can be used to specify any layer 4 protocol we want to use, so for example `IPPROTO_IP` is not supported as the `protocol` since it is a dummy layer 3 protocol. Note that we can only send and receive a **single protocol** using this, it is not possible to receive data for all protocols with a RAW-socket [1]. We also have to be aware that the kernel still also receives all protocol data and will act accordingly [1]. If we do not want the kernel to also react to the packets we have to take additional steps to prevent this.

Unix RAW-sockets work as described in Linux distributions. For Unix based systems there are some limitations we have to observe.

- The Write-commands work for layer 3 and 4 work for Unix as well.
- The Read-commands for layer 3 and 4 only work with some limitations. TCP and UDP packets cannot be received with RAW-sockets, they are only delivered to the kernel [6]. Copies of ICMP packets are delivered to the RAW-socket as well as to the kernel [6]. All IGMP packets are delivered to the RAW-socket, as well as any other protocol packets [6].
- Unix does not support the PACKET-socket as described for layer 2 later. Instead Unix has the Berkeley Packet Filter (BPF) interface that can be used for layer 2 communication [2]. We will give a short introduction in a later section.
- Some of the functions that are used in this section are requiring additional headers to work under Unix. The additional headers are noted in the last section where the functions were introduced.
- Other than these limitations we can work with the RAW-socket under Unix in the same way as in Linux [8].

Read The structure of the complete Read operation for a basic layer 4 RAW-socket is shown in figure 3 on the next page.

- We first have to create the socket as it was discussed in section 3.2.1 on page 15. Like mentioned, if we only want to create layer 4 headers later we can choose any protocol family we want to receive data for except the `IPPROTO_RAW` [1]. Also even if we just read data to the buffer it makes sense to also reserve buffer space and initialize it with zero to have defined data in the buffer.
- As an optional step we can define the interface we want to receive the network data from by defining the source address using the `bind()` function.
- Additionally we can, as an option, use the `connect()` function to define a standard source we want to receive data from or we want to connect to in case of connection-oriented protocols [8].
- Then we can read binary data from the socket with one of the function described in section 3.2.7 on page 19. We can either read data using the `read()` or `recv()` functions. As an alternative we can use the `recvfrom()` function which also returns the source address of the packet we did receive. It is recommended to check the return value to know how many bytes were received to be able to react accordingly.
- Optionally we can `close()` the socket now if we received all the data we want to receive.
- Then we can type-cast the binary data onto our headers as it is described in section . We have to be aware that the Buffer of a RAW-socket includes the headers as well as the data . The included headers start at the Network layer (OSI layer 3) [1]. Since we do not want to use the information of the Network layer at this point, we have to add the length of its header to the data pointer to get the Transport layer (OSI layer 4) header information [5]. When doing this we have to be aware that there can be variable network layer headers as described in section 3.1.4 on page 14. Depending on the protocol and possibly additional variable headers we can compute the start positions of the next header (or headers) and type-cast them into structures.

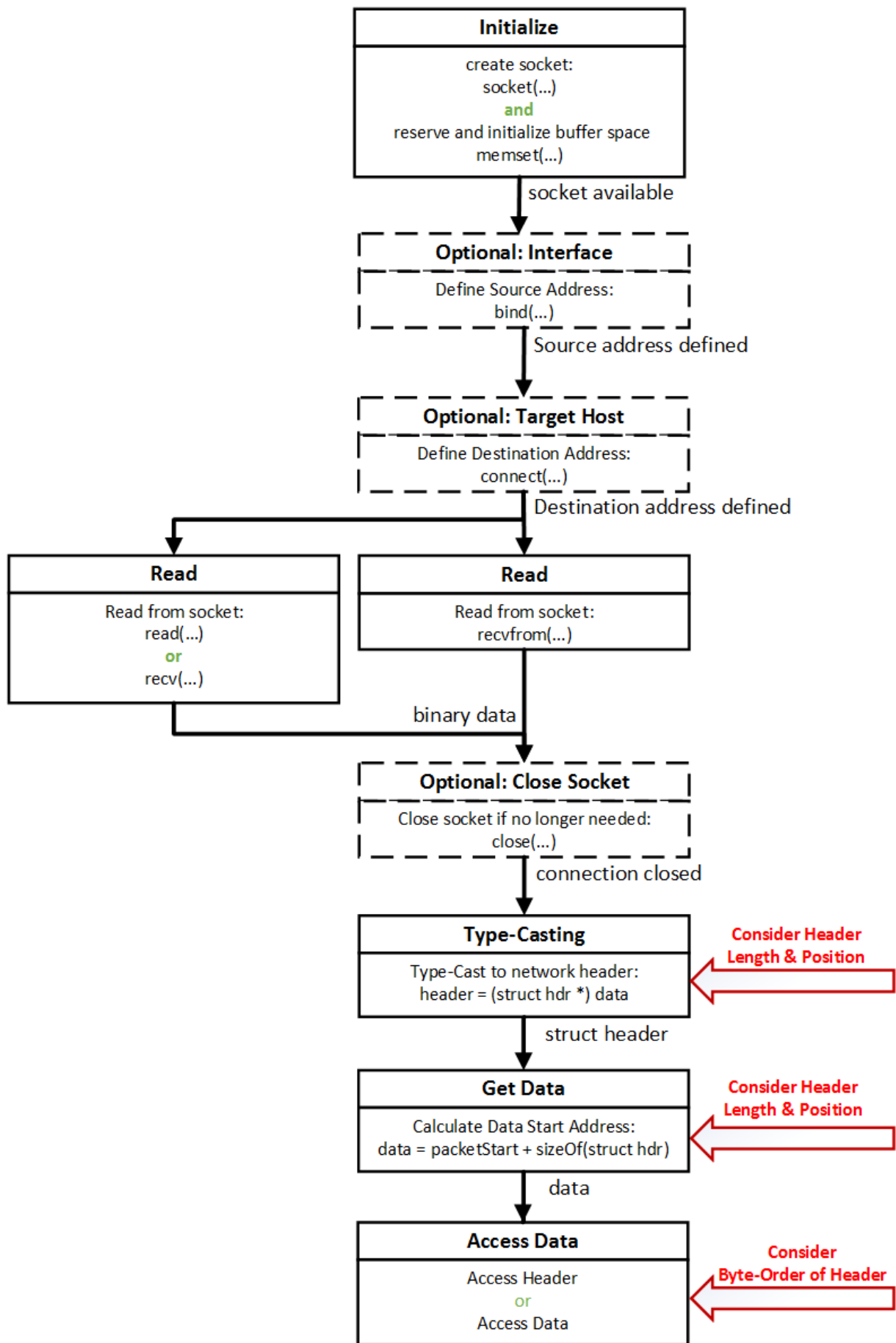


Figure 3: Structure of a RAW-socket layer 4 Read Operation

- After all headers are handled the remainder of the packet is the user data and can also be extracted.
- Then we can access the extracted headers or data. When accessing the headers we have to observe the byte-order, as discussed in the section 3.1.1 on page 13.

As the general programming paradigms suggest, the user should handle possible errors that might occur in the steps. Possible error values are presented in the introductions of the available functions. This is even more important when working with RAW-sockets, since working with them is very error prone.

Write When we use the layer 4 Write access, the operating system will take over most function required to send the packet. We only have to start to create our packet at layer 4 and supply the required destination address to the API to send the packet. All layer 3 headers will be filled in by the operating system accordingly [1].

- We first have to create the socket as it was discussed in section 3.2.1 on page 15. Like mentioned if we only want to create layer 4 headers, we can choose any protocol family except the `IPPROTO_RAW`[1]. Also we reserve buffer space for the packet we want to create and initialize it with zero to have defined data in the buffer [5].
- As an optional step we can define the interface we want to use to send the network data from by defining the source address using the `bind()` function.
- Then we can type-cast the buffer to our layer 4 headers like it is described in section 3.3.7 on page 44. We have to take into consideration variable header lengths if we want to use optional fields and compute the start of the next header accordingly as described in section 3.1.4 on page 14. The remaining buffer can be used for the user data we might want to pass along.
- As the next step we set the data of our header as required. We have to consider the byte-order in case of fields that span multiple bytes and change the byte-order according to the requirements of the protocol we want to use [5]. Additional information about this can be found in section 3.1.1 on page 13.
- If the protocol does use a checksum to protect the header from changes we also have to compute the checksum for the layer 4 ourselves [5]. Only for the layers below, the operating system will take care of computing the checksum. We have to see in the documentation which kind of checksum is required by the protocol we want to use. Additional information about checksum computation can be found in the section 3.1.2 on page 13.
- In the next step we have the choice to use one of two possibilities. We can use the `sendto()` function, which allows us to pass a destination address along with the data. Optionally we can as an option use the `connect()` function to define a standard source we want to send data to or want to connect to in case of connection-oriented protocols. If we define the destination with `connect()` the destination address of `sendto()` might be left empty, then the already defined address is used [8]. The other possibility is to use `write()` or `send()` which both require the socket to be in a connected state, so we have to call `connect()` first to use one of these functions [8]. All functions have the same result, the data is being sent. The result these functions return is the number of bytes that were sent. It is recommended to check the return value and match it with the length of the data that should have been sent to make sure the function was not interrupted during the call and all data has been sent as intended [8].
- Optionally we can `close()` the socket now if we already sent all the data we want to transmit.

3.2.14 Layer 3

Using RAW-sockets to also access the layer 3 of a network packet is very similar to the layer 4 access we discussed in the previous section. The most significant differences are that we need to set a different option to get write access to the complete layer 3 header. We also now have to take the structure and lengths of multiple headers into consideration. For example the basic call to create a raw-socket for IP could look like this [5]:

```
sockfd = socket(AF_INET, SOCK_RAW, IPPROTO_RAW)
```

In this example the `AF_INET` specifies IP as the protocol, the `SOCK_RAW` defines the socket type as a RAW-Socket and the `IPPROTO_RAW` specifies that we are interested in sending any type of protocol with this socket. Note that we can only send any protocol with this, it is not possible to receive any data with this socket [1]. If we want to still receive packets we could specify a single protocol in the last parameter of the `socket()` call and then have to set the `IP_HDRINCL` with the `setsockopt()` function to be able to write and receive with the same socket. We also have to be aware that the kernel still also receives all protocol data and will act accordingly [1]. If we do not want the kernel to also react to the packets we have to take additional steps to prevent this.

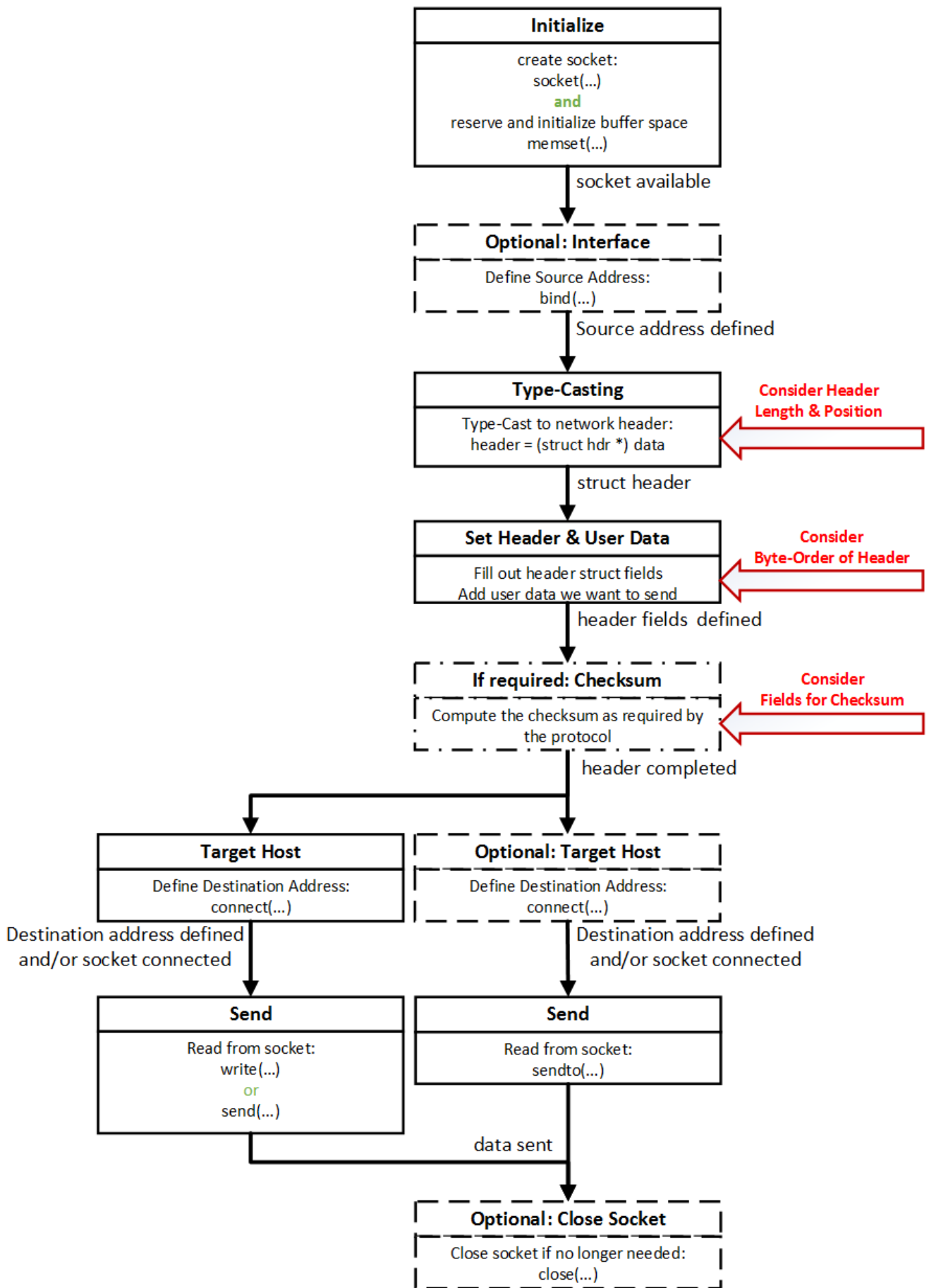


Figure 4: Structure of a RAW-socket layer 4 Write Operation

Read For the Read operation the structure and commands are the same as for the layer 4 read access. As already discussed we simply ignored the layer 3 information before since we did not want to use it. Nevertheless the same command already gave us the possibility to read-out the header information. Also if we also want to have write access to the layer 3 headers for the write operation later it is a possible simplification to create a socket and then set the `IP_HDRINCL` flag as already discussed. This is not necessary for the Read operation we want to discuss first only for the Write later. Details on this option will follow in the layer 3 Write section. It should be noted again that receiving data of all IP-protocols is not possible with RAW-sockets, only the data for a single protocol can be received at a time [1]. A possible solution for this are the PACKET-sockets discussed in the next section.

The structure of the complete Read operation for a basic layer 3 RAW-socket is shown in figure 5 on the next page.

As already mentioned there are only slight differences due to the handling of multiple layers in comparison to the layer 4 read.

- We first have to create the socket as it was discussed in section 3.2.1 on page 15. Like mentioned if we want to create layer 3 headers we can additionally set the `IP_HDRINCL` flag with the `setsockopt()` function. Also even if we just read data to the buffer it makes sense to also reserve buffer space and initialize it with zero to have defined data in the buffer [5].
- As an optional step we can define the interface we want to receive the network data from by defining the source address using the `bind()` function.
- Additionally we can as an option use the `connect()` function to define a standard source we want to receive data from or we want to connect to in case of connection-oriented protocols.
- Then we can read binary data from the socket with one of the function described in section 3.2.7 on page 19. We can either read data using the `read()` or `recv()` functions. As an alternative we can use the `recvfrom()` function which also returns the source address of the packet we received. It is recommended to check the return value to know how many bytes were received to be able to react accordingly [8].
- Optionally we can `close()` the socket now, if we received all the data we want to receive.
- Then we can type-cast the binary data into our headers like it is described in section 3.3.7 on page 44. We have to be aware that the buffer of a RAW-socket includes the headers as well as the data [1]. The included headers start at the network layer (OSI layer 3). Since we want to use the information of the network layer at this point, we can typecast the header onto a struct to get access to all the fields. We still have to add the length of its header to the data pointer to get the Transport layer (OSI layer 4) header information. When doing this we have to be aware that there can be variable network layer headers as described in section 3.1.4 on page 14. Also depending on which functions we used for reading data from the socket, we might want to filter from which source and/or for which target address we want to receive data from. Depending on the protocol and possibly additional variable headers we can compute the start positions of the next header (or headers) and type-cast them into structs.
- After all headers are handled the remainder of the packet is the user data that can also be extracted.
- Then we can access the extracted headers or data. When accessing the headers we have to observe the byte-order, as discussed in the section 3.1.1 on page 13.

Write Like already noted, to get write access to the layer 3, we have to either create an `IPPROTO_RAW` socket when choosing our protocol for the `socket()` function or set the `IP_HDRINCL` with the `setsockopt()` to deactivate automatic IP-header generation by the operating system and be able to manually fill the fields [5]. It should be noted here that `IPPROTO_RAW` socket does not include the `IP_HDRINCL` flag for all operating systems since this is not defined in the POSIX-Standard [11]. So if you want to be sure that the code works regardless of the distribution used, it is recommended to set the **protocol** option to the protocol that can be found in the layer 3 protocol header and use `setsockopt()` to tell the OS that we want to generate the header manually [11].

The Structure of the complete Write operation for a basic layer 3 RAW-socket is shown in figure 6 on page 34.

- We first have to create the socket as it was discussed in section 3.2.1 on page 15. Like mentioned we have to use either `IPPROTO_RAW` when choosing our protocol for the `socket()` function, or set the `IP_HDRINCL` with the `setsockopt()` function to deactivate automatic IP-header generation as already discussed [5]. Also we reserve buffer space and initialize it with zero to have defined data in the buffer [5].

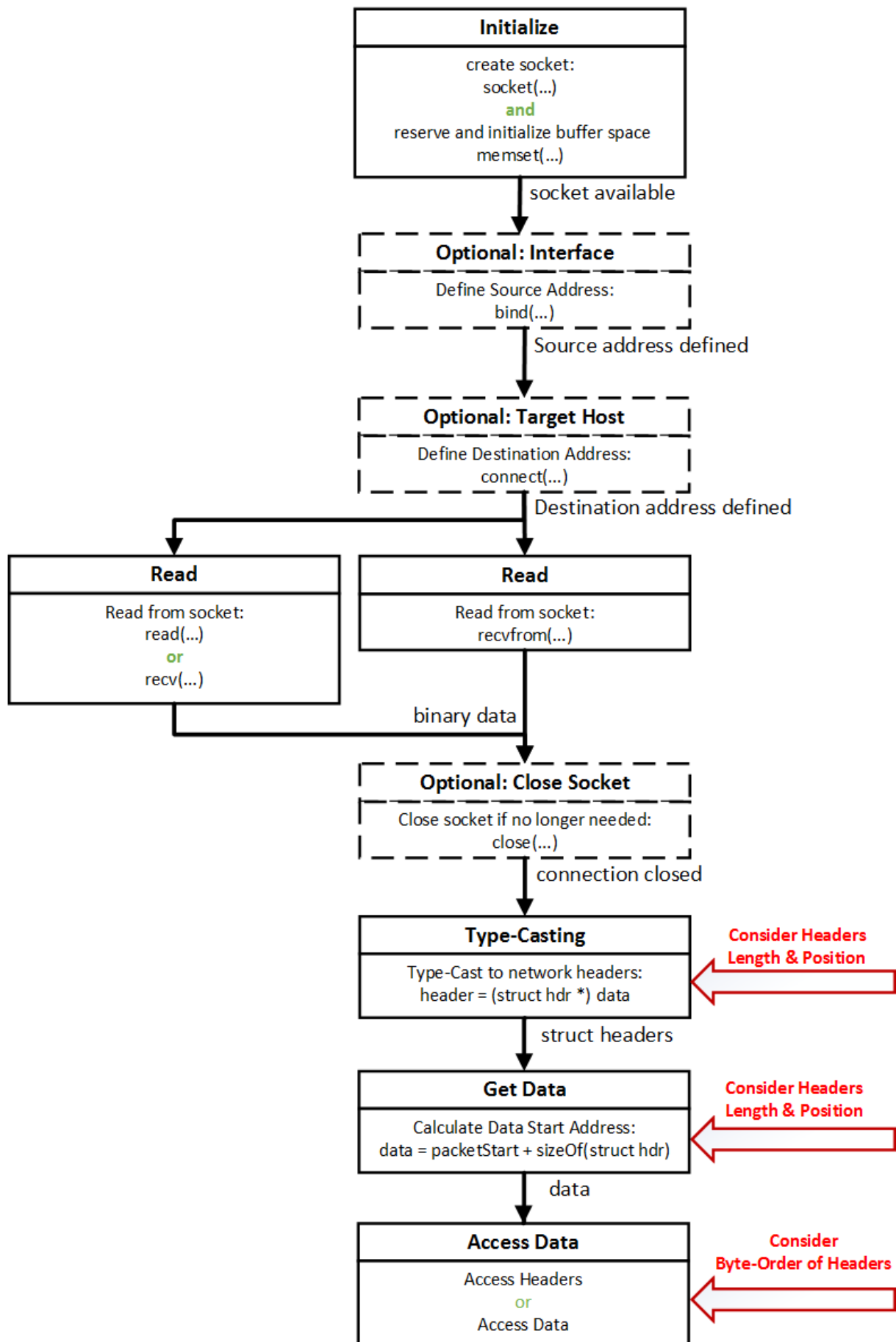


Figure 5: Structure of a RAW-socket layer 3 Read Operation

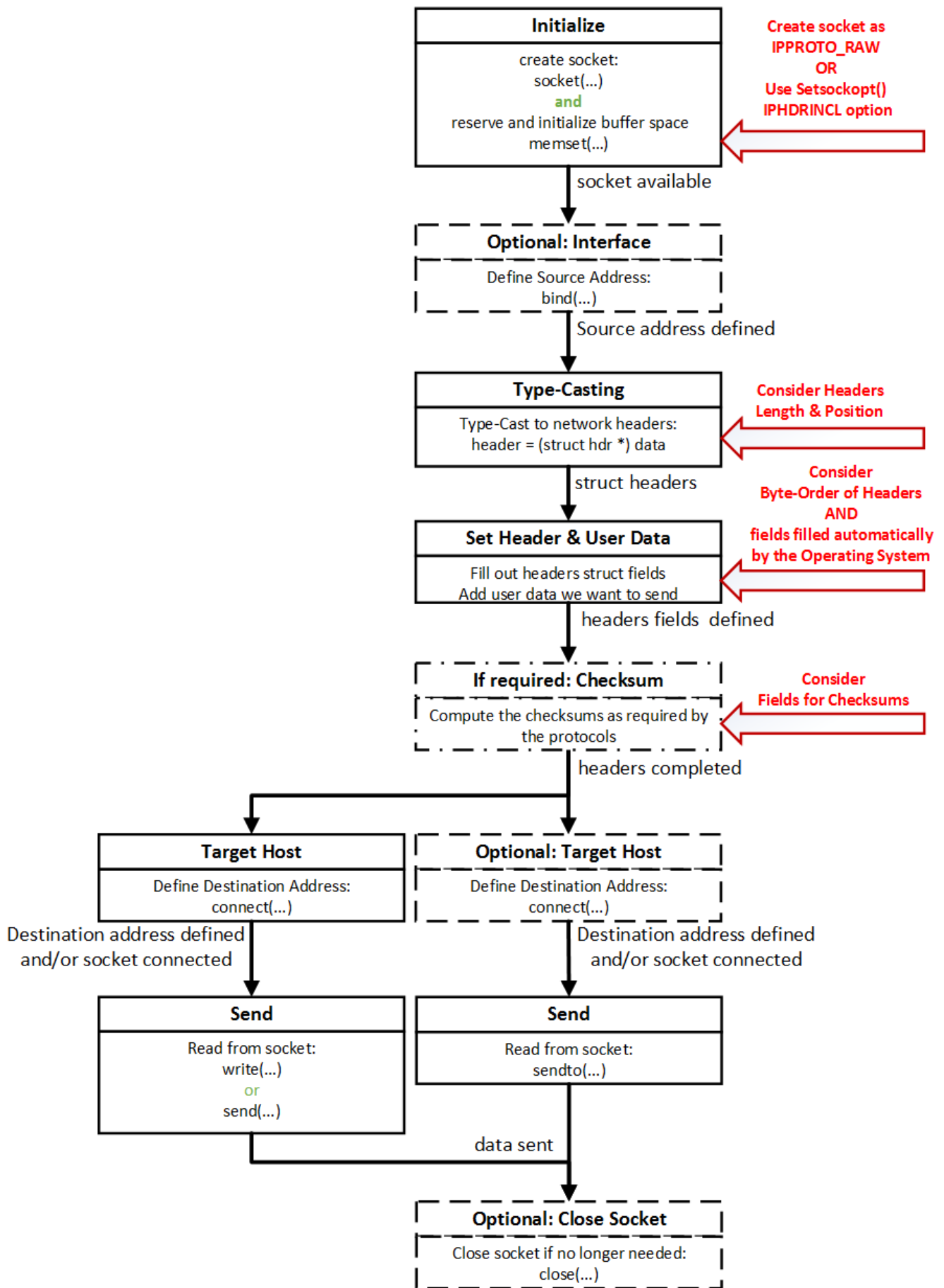


Figure 6: Structure of a RAW-socket layer 3 Write Operation

- As an optional step we can define the interface we want to use to send the network data with by defining the source address using the `bind()` function.
- Then we can type-cast the buffer to our headers like it is described in section 3.3.7 on page 44. The included headers start at the Network layer (OSI layer 3). We have to take into consideration variable header lengths, if we want to optional fields and compute the start of the next header accordingly as described in section 3.1.4 on page 14. The remaining buffer is reserved for the user data we might want to pass along.
- As the next step we set the data of our headers as required. However, we do not need to fill the checksum and total length of the IP-packets, since the operating system will fill in the values automatically [1]. Optionally we can also fill in zero in the fields source address and packet ID [1]. If we do this then the operating system also will fill in the values automatically. When we fill in the values of the header fields we have to consider the byte-order in case of fields that span multiple bytes and change the byte-order according to the requirements of the protocol we want to use. Additional information about this can be found in section 3.1.1 on page 13.
- If the protocol does use a checksum to protect the header from changes, we also have to compute the checksum for the layer 4 ourselves [5]. Only for the layers below, the operating system will take care of computing them. For the IP-Protocol for example the Operating System will always fill in the value for the checksum. We have to see in the documentation which kind of checksum is required by the protocol we want to use. Additional information about computing the checksum can be found in the section 3.1.2 on page 13.
- In the next step we have the choice to use one of two possibilities. We can use the `sendto()` function, which allows us to pass a destination address along with the data. Optionally we can as an option use the `connect()` function to define a standard source we want to send data to, or want to connect to in case of connection-oriented protocols. If we define the destination with `connect()` the destination address of `sendto()` might be left empty, then the already defined address is used. The other possibility is to use `write()` or `send()` which both require the socket to be in a connected state, so we have to call `connect()` first to use one of these functions [8]. All functions have the same result, the data is being sent. It is recommended to check the return value and match it with the length of the data that should have been sent to make sure the function was not interrupted during the call and all data has been sent as intended [8].
- Optionally we can `close()` the socket now, if we already sent all the data we want to transmit.

3.2.15 Layer 2 - PACKET-sockets

In addition to the layers we discussed so far, it is also possible under Linux to access the Data-Link-layer with RAW-sockets. In earlier versions there was a special type of sockets for this, the **PACKET_socket** [5]. Nowadays the use of packet sockets is no longer recommended, instead we can use a **normal RAW_socket** with a different protocol supplied to the `socket()` function, like shown in section 3.2.1 on page 15 [5]. The possible socket types we can use are the `SOCK_RAW` to get access to the whole packet including the Data-Link-layer header or the `SOCK_DGRAM` which operates on the Data-Link-layer packet without the headers [12]. The protocols for Ethernet we can use in Linux are shown in the table 48 on page 81. One common option might be the `ETH_P_ALL` protocol constant which results in all incoming and outgoing Ethernet packets to be accessed with the RAW-socket. A typical call looks like this, note the `htons` function used to change the byte-order of the supplied protocol constant since it is a multi-byte constant [5]:

```
s = socket(AF_PACKET, SOCK_RAW, htons(ETH_P_ALL));
```

With this all headers down to the Data-Link-layer can now be accessed and all data received on all interfaces can be received. However as already mentioned, parts of the Data-Link-layer header might not be accessible, for example for the Ethernet frames the CRC (Cyclic Redundancy Check) checksum, the Preamble and the Start-of-Frame Delimiter are not accessible since these fields are more part of the Physical-layer and therefore handled by the network driver [5]. Also with PACKET-socket the `connect()` function has lost its purpose, only the `bind()` function still makes sense to define the network interface [5]. We have to be aware that the kernel still also receives all protocol data and will act accordingly [1].

Promiscuous Mode The promiscuous mode is a special mode of a network card in which all network traffic of an interface is received, in contrast to the usual case where only the traffic addressed to the interface is received [5]. PACKET-sockets work similar to the promiscuous mode, but receive all traffic for all interfaces of the host, but can be configured to only receive the traffic on a single interface by using the `bind()` function [5]. If it is necessary to activate the promiscuous mode for some reason it can be done by setting the `IFF_PROMISC` using the `ioctl()` systemcall [5].

MAC-Address Since we now work on the layer 2 we also have to work with the addresses of the layer, for example the MAC-addresses used for example in Ethernet. The MAC-address of the interface we want to use as a source can be accessed with a call of the `ioctl` function. We demonstrate the call in the following listing. Note that the listing does not include any protection against errors.

Listing 7: mac.cpp

```
1 struct ifreq ifr; // struct for the interface information
2 char * eth = "wlan0"; // name of the interface
3
4 // copy the name to the struct
5 memcpy(ifr.ifr_name, if_name, IFNAMSIZ);
6
7 //create a socket descriptor
8 int sockfd = socket( AF_PACKET , SOCK_RAW , htons(ETH_P_ALL)) ;
9
10 //get the information
11 ioctl(sockfd, SIOCGIFHWADDR, &ifr)
12
13 //extract the mac
14 const unsigned char* mac = (unsigned char*) ifr.ifr_hwaddr.sa_data;
```

Read Except for the creation of the packet socket, the missing `connect()` statements and the handling of more layers, the layer 2 read is not much different from the other read operations and share the same basic structure. The Structure of the complete Read operation for a basic layer 2 PACKET-socket is shown in figure 7 on the next page.

For a read of the layer 2 the general structure we discussed until now is the same, only some steps are different due to different initialization we need and the handling of Data-Link-layer data:

- We first have to create the socket as it was discussed in section 3.2.1 on page 15. As already discussed the possible socket types we can use are the `SOCK_RAW` or the `SOCK_DGRAM`. The protocols for Ethernet we can

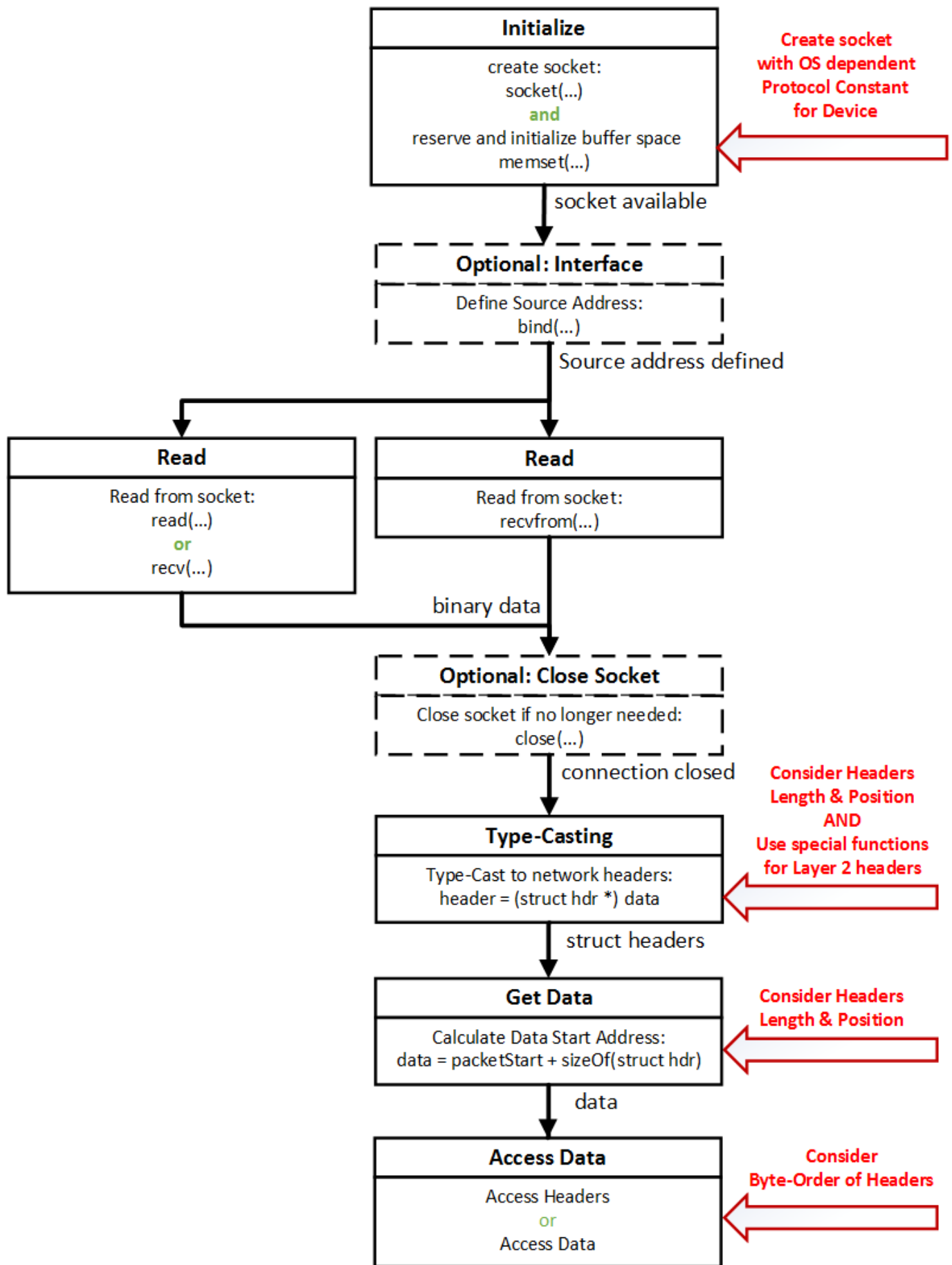


Figure 7: Structure of a PACKET-socket layer 2 Read Operation

use in Linux are shown in the table 48 on page 81. Since they could be multi-byte constants we have to use functions to change the byte-order as discussed in the section 3.1.1 on page 13. Also even if we just read data to the buffer it makes sense to also reserve buffer space and initialize it with zero to have defined data in the buffer [5].

- As an optional step we can define the interface we want to receive the network data from by defining the source address using the `bind()` function.
- Then we can read binary data from the socket with one of the function described in section 3.2.7 on page 19. We can either read data using the `read()`, `recv()` functions, or we can use the `recvfrom()` function which also returns the source address of the packet we did receive. It is recommended to check the return value to know how many bytes were received to be able to adapt accordingly [8].
- Optionally we can `close()` the socket now, if we received all the data we want to receive.
- Then we can type-cast the binary data into our headers like it is described in section . Since we did start at the Data-Link-layer, there could be different layer 3 protocols included in the frame. This is why we have to access the correct field in the header to find out what is encapsulated in the frame and which header structure we have to handle. When doing this we have to be aware that there can be variable network layer headers as described in section 3.1.4 on page 14. Also depending on which functions we used for reading data from the socket we need to filter from which source and/or for which target address we want to receive data from. Depending on the protocol and possibly additional variable headers we can compute the start positions of the next header (or headers) and type-cast them into structs.
- After all headers are handled the remainder of the packet is the user data that can also be extracted.
- Then we can access the extracted headers or data. When accessing the headers we have to observe the byte-order, as discussed in the section 3.1.1 on page 13.

Write The Write on a PACKET-socket differs from the Write on a RAW-socket in some important parts. The IP-header fields that were filled by the operating system when using a RAW-socket have to be filled in manually when using a PACKET-socket. The only layer where the operating system still helps the user is on the Data Link Layer, where some fields are automatically generated [5]. That is the case for example when creating Ethernet frames, then the Padding field of the frame, the Preamble, the Start-of-Frame Delimiter and the CRC checksum [5]. These fields are still generated automatically by the network driver.

The Structure of the complete Write operation for a basic layer 2 PACKET-socket is shown in figure 8 on the following page.

- We first have to create the socket as it was discussed in section 3.2.1 on page 15. As already discussed the possible socket types we can use are the `SOCK_RAW` or the `SOCK_DGRAM`. The protocols for Ethernet we can use in Linux are shown in the table 48 on page 81. Since they could be multi-byte constants we have to use functions to change the byte-order as discussed in the section 3.1.1 on page 13. Also we reserve buffer space and initialize it with zero to have defined data in the buffer.
- Even if packets sent over PACKET-sockets have to be generated completely manual it is still required for the send functions to define a socket-address. The necessary struct `sockaddr_ll` can be found in the header `<linux/if_packet.h>`. The only part of the struct we have to fill in is the network interface `sll_ifindex`[5]. Now to find the network interface, in Linux usually denoted as f.e. `eth1` we can use the so called *Netdevice-Interface* [10]. It describes the `ioctl()`-system calls that can be used to access different properties of network interfaces in Linux, for example the callcode `SIOCGIFINDEX` has to be used to get the desired information. The access can be done on all socket types with the data structure `ifreq` that is defined in `<linux/net/if.h>`. The complete call we have to use looks like this [5]:

```
ioctl(socketfd, SIOCGIFINDEX, &ifr)
```

The `socketfd` again is our socket descriptor, the `SIOCGIFINDEX` is the parameter to retrieve the information we need and `ifr` is the `ifreq` that after the call has the information we need. With the information we can later call the `bind()` or `sendto()`.

- Then we can type-cast the buffer to our headers like it is described in section . The included headers start at the Data-Link-layer (OSI layer 2). It is necessary to include or define additional headers to have structs for the Data-Link-layer header. We have to take into consideration variable header lengths, if we want to optional fields and compute the start of the next header accordingly as described in section 3.1.4 on page 14. The remaining buffer is reserved for the user data we might want to pass along.

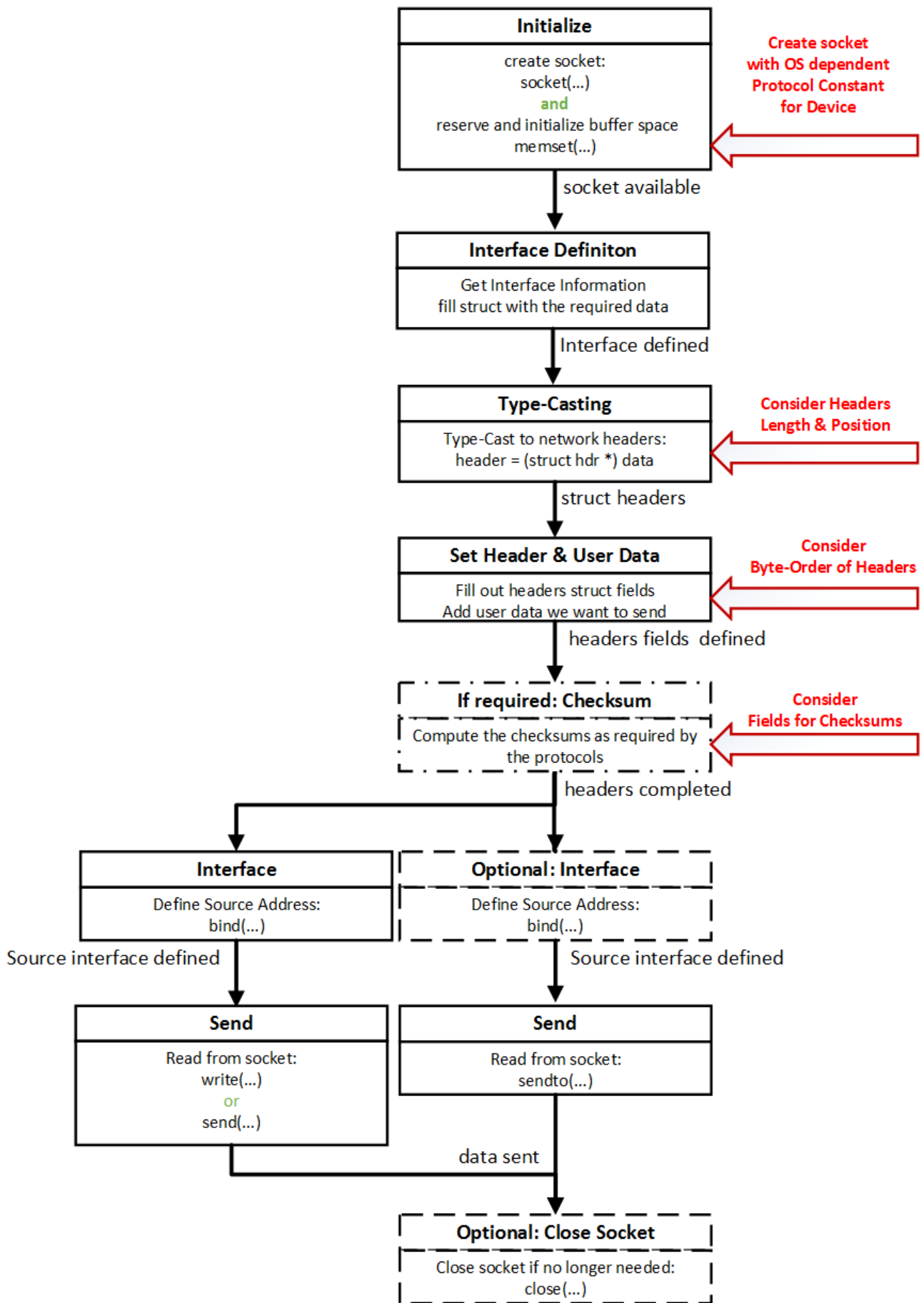


Figure 8: Structure of a PACKET-socket layer 2 Write Operation

- As the next step we set the data of our headers as required. We have to fill in all fields since the operating system does not automatically create fields for PACKET-sockets [5]. When we fill in the values of the header fields have to consider the byte-order in case of fields that span multiple bytes and change the byte-order according to the requirements of the protocol we want to use. Additional information about this can be found in section 3.1.1 on page 13.
- If the protocol does use a checksum to protect the header from changes, we have to compute them and set the fields accordingly [5]. We have to see in the documentation which kind of checksum is required by the protocol we want to use. Additional information about computing a checksum can be found in the section 3.1.2 on page 13.
- In the next step we have the choice to use one of two possibilities. We can use the `sendto()` function, which allows us to pass the socket address of the interface we created along with the data. Optionally we can as an option use the `bind()` function to define a standard socket address for the interface we want to send data with. If we define the destination with `bind()` the destination address of `sendto()` might be left empty, then the already defined address is used [8]. The other possibility is to use `write()` or `send()` which both require a call of `bind()` first. All functions have the same result, the data is being sent. It is recommended to check the return value and match it with the length of the data that should have been sent to make sure the function was not interrupted during the call and all data has been sent as intended [8].
- Optionally we can `close()` the socket now, if we already sent all the data we want to transmit.

3.3 Berkeley Packet Filter (BPF)

The Berkeley Packet Filter (BPF) provides us with the same functionality we have introduced for the PACKET-socket in Linux, it provides full access to the data link layer for read and write [2]. So as it is the case for the PACKET-socket we can receive any packet received by any interface regardless if it is meant for our host or not [2]. Additionally to these functions the BPF also includes an advanced filtering mechanism that we can use to filter incoming packets according to different criteria [2]. We will only give a short overview over this function since our main focus is not to give an introduction on how to program a network sniffer, but how to send and receive data including the protocol headers.

To use the BPF we have to use a distribution that has the device `bpf` in its kernel [7]. Only then can we create a BPF device. The BPF appears as a special character device `/dev/bpf`. To open it we have to use the well known `open()` function which is defined in `<fcntl.h>` and looks like this:

```
int open(const char * pathname, int flags)
```

We have to set the following parameters:

- **pathname** takes the path of a file. In our case the path has to be `/dev/bpf n` where n is the number of the device depending on how many other programs use a BPF [7].
- **flags** takes a flag as the name implies, the most useful for us is the `O_RDWR` flag that gives us read and write access.
- After we have done this we have to associate the BPF device to a network interface [7]. This works similar to the `bind()` function used for the RAW- or PACKET-socket, but is realized with a `ioctl()` call with the flag `BIOCSETIF`. The completed command looks like this, note the `ifreq` structure:

Listing 8: bpf.cpp

```
1 // define the name of the interface
2 const char* interface = "fxp0";
3
4 // create the struct for the interface
5 struct ifreq bound_if;
6
7 // copy the name into the struct
8 strcpy(bound_if.ifr_name, interface);
9
10 // associate the bpf to the interface
11 ioctl( bpf, BIOCSETIF, &bound_if )
```

Afterwards we just have to set the BPF to immediate mode, since a `read()` would otherwise block until the kernel buffer becomes full or a timeout occurs [2]. We also have to get the buffer size, since the BPF might return us more than one packet at a time [7]. The code for this is shown in the listing below:

Listing 9: bpf2 cpp

```

1 // initialize buffer length
2 int buf_len = 1;
3
4 // activate immediate mode
5 ioctl( bpf, BIOCIMMEDIATE, &buf_len )
6
7 // request buffer length
8 ioctl( bpf, BIOCGLEN, &buf_len )

```

After this we can read and write to the BPF like we did it before. Note that we only can use the `read()` and `write()` functions since the BPF is handled like a file. Before we introduce these operations we want to show some BPF specialties and give an overview over some preferences that can be set for the BPD subsystem.

3.3.1 BPF Header

The BPF device adds an additional header with information about the received packet [2]. One of the headers shown in the following listing is appended to each packet, `bpf_xhdr` is used by default, `bpf_hdr` is used only when the timestamp format is set to certain values using the `ioctl()` function [2].

Listing 10: bpf structs cpp

```

1 struct bpf_xhdr {
2     struct bpf_ts    bh_tstamp;    /* time stamp */
3     uint32_t         bh_caplen;    /* length of captured portion */
4     uint32_t         bh_datalen;    /* original length of packet */
5     u_short          bh_hdrlen;    /* length of bpf header (this struct
6                                     plus alignment padding) */
7 };
8
9 struct bpf_hdr {
10     struct timeval   bh_tstamp;    /* time stamp */
11     uint32_t         bh_caplen;    /* length of captured portion */
12     uint32_t         bh_datalen;    /* original length of packet */
13     u_short          bh_hdrlen;    /* length of bpf header (this struct
14                                     plus alignment padding) */
15 };

```

The `bh_hdrlen` is used for padding between the header and the link layer protocol, additionally each packet is padded so that it starts on a word [2]. A macro `BPF_WORDALIGN()` is defined in the header that allows us to compute the header-position correctly when we compute the header-positions.

3.3.2 Buffer Modes

BPF devices can deliver packet data to the applications in two different modes, **Buffered read mode** and **Zero-copy mode** [2]. The mode can be set by using the `ioctl()` function using the `BIOCSETBUFMODE` flag. The default mode is the **Buffered read mode** which can also be set using the `BPF_BUFMODE_BUFFER` flag. In this mode the packet data can be accessed by explicitly call the `read()` function. A fixed buffer size is used for all internal buffers as well as for the `read()` function which can be queried using the `BIOCGLEN` flag for the `ioctl()` call and set by using the `BIOCSLEN` flag in the `ioctl()` call [2]. Note that packets longer than this buffer size are truncated. The other mode is the **Zero-copy buffer mode** that can be set using the `ioctl()` call and the `BPF_BUFMODE_ZEROCOPY` flag. In this mode the user process registers two equal sized buffers using the `BIOSETZBUF` flag with the `ioctl()` function in which the packet data is directly stored [2]. The user process than has to use atomic operations to check if it can read data from a buffer and than return it to the kernel for storing the next data as fast as possible [2]. See the Unix MAN-pages for detailed information how to use this mode.

3.3.3 IOCTLS

The following flags can be used to change the behavior of the BPF device. All constants are defined in the `<bpf.h>` header, except for the `BIOGETIF` and `BIOSETIF` which also require the `<sys/socket.h>` and `<net/if.h>` headers. They can be set using the `ioctl()` function for any open BPF file, with the type indicated as the third argument of the function. See table 52 in the appendix.

3.3.4 SYSCALL Variables

A set of syscall-Variables for controlling the behavior of the BPF subsystem exist [2]:

Variable	Description
<code>net.bpf.optimize_writers</code>	Turning this option on makes new BPF users to be attached to write-only interface list until program explicitly specifies read filter via <code>pcap_set_filter()</code> .
<code>net.bpf.stats</code>	Binary interface for retrieving general statistics.
<code>net.bpf.zerocopy_enable</code>	Permits zero-copy to be used with net BPF readers.
<code>net.bpf.maxinsns</code>	Maximum number of instructions that BPF program can contain.
<code>net.bpf.maxbufsize</code>	Maximum buffer size to allocate for packets buffer.
<code>net.bpf.bufsize</code>	Default buffer size to allocate for packets buffer.

Table 20: `ioctl()` flags defined in `<bpf.h>` [2]

3.3.5 Filter Maschine

An additional capability we have is the Filter machine we want to just show here. It allows to filter the incoming packets for specific packet types. To do this we can define an array of instructions which is executed for every received packet and which consists of an accumulator, index register, scratch memory store, and implicit program counter [2]. The following listing shows the struct:

Listing 11: bpf structs filter cpp

```
1 struct bpf_insn {
2     u_short code;
3     u_char  jt;
4     u_char  jf;
5     u_long  k;
6 };
```

The `k` field is used in different ways by different instructions, and the `jt` and `jf` fields are used as offsets by the branch instructions [2]. There are eight classes of instructions like for example jump and value copy [2]. Various other mode and operator bits are inserted into the class to give the actual instructions [2]. Further information about the filters can be found in the Unix manpages.

3.3.6 Read

The read of a packet using an BPF device is similar to the same operation in a packet socket. The basic structure is shown in figure 9 on the next page.

For a read operation with the BPF we have to do the following:

- We first have to create the BPF device as discussed in the beginning of this section. Also even if we just read data to the buffer it makes sense to also reserve buffer space and initialize it with zero to have defined data in the buffer.
- We then have to attach the BPF device to a interface as discussed in the beginning of this chapter to be able to receive frames [7]. We use a `ioctl()` call to achieve this.
- Then we can set options, like the `BIOCIMMEDIATE` mode which was discussed earlier, using `ioctl()` calls [7]. We also could set a filter to only receive certain packets we are interested in [2].
- As the next step we retrieve the buffer length from the kernel to know how many bytes are received per packet [7]. That is necessary since the received packet also includes a BPF header containing information about the packet received as discussed in the earlier sections.
- Then we can read binary data from the socket with the read function described in section 3.2.7 on page 19. It is recommended to check the return value to know how many bytes were received to be able to adapt accordingly [7]. If we want to read multiple packets from the BPF device we have to use the `BPF_WORDALIGN()`-macro to align the pointer for the next header correctly [7].
- Optionally we can `close()` the BPF device now, if we received all the data we want to receive.

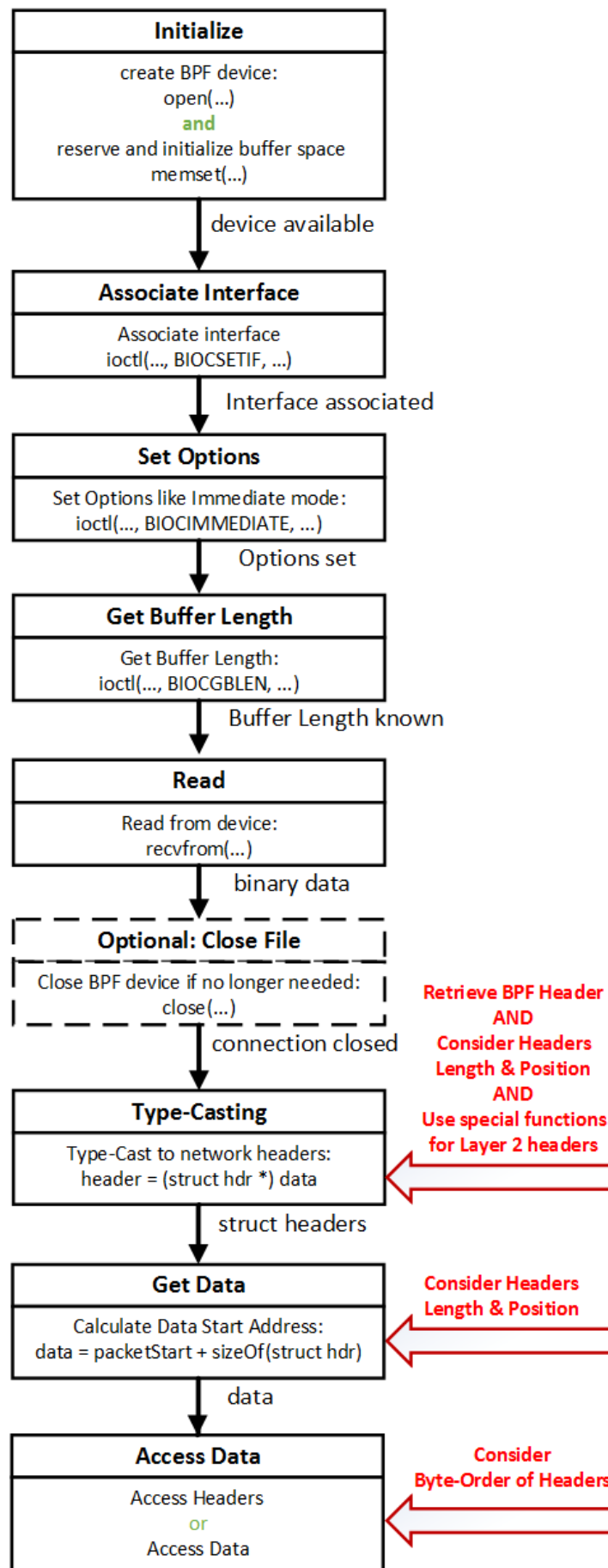


Figure 9: Structure of a BPF device layer 2 Read Operation

- Then we can type-cast the binary data into our headers like it is described in section . We have to take into consideration that there is an additional BPF header attached to the received packet which has to be retrieved first [7]. Since we did start at the Data-Link-layer, there also could be different layer 3 protocols included in the frame. This is why we have to access the correct field in the header to find out what is encapsulated in the frame and which header structure we have to handle. When doing this we have to be aware that there can be variable network layer headers as described in section 3.1.4 on page 14. Also depending on which functions we used for reading data from the socket we need to filter from which source and/or for which target address we want to receive data from. Depending on the protocol and possibly additional variable headers we can compute the start positions of the next header (or headers) and type-cast them into structs.
- After all headers are handled the remainder of the packet is the user data that can also be extracted.
- Then we can access the extracted headers or data. When accessing the headers we have to observe the byte-order, as discussed in the section 3.1.1 on page 13.

3.3.7 Write

As figure 9 on the previous page shows the basic structure, the write operation of a packet using an BPF device is also similar to the same operation in a packet socket.

- We first have to create the BPF device as discussed in the beginning of this section [7]. Also even if we just read data to the buffer it makes sense to also reserve buffer space and initialize it with zero to have defined data in the buffer.
- We then have to attach the BPF device to a interface as discussed in the beginning of this chapter to be able to receive frames [7]. We use a `ioctl()` call to achieve this.
- Then we can set options, like the `BIOCIMMEDIATE` mode which was discussed earlier, using `ioctl()` calls. It also is possible to set the `BIOCGHRCMPLT` option, which select if the operating system should auto fill in the data link headers[2].
- Then we can type-cast the buffer to our headers like it is described in section . The included headers start at the Data-Link-layer (OSI layer 2). It is necessary to include or define additional headers to have structs for the Data-Link-layer header. We have to take into consideration variable header lengths if we want to optional fields and compute the start of the next header accordingly as described in section 3.1.4 on page 14. The remaining buffer is reserved for the user data we might want to pass along.
- As the next step we set the data of our headers as required. We can have the kernel auto fill the data link layer addresses if we used the `BIOCGHRCMPLT` option [2]. When we fill in the values of the header fields have to consider the byte-order in case of fields that span multiple bytes and change the byte-order according to the requirements of the protocol we want to use. Additional information about this can be found in section 3.1.1 on page 13.
- If the protocol does use a checksum to protect the header from changes, we have to compute them and set the fields accordingly [5]. We have to see in the documentation which kind of checksum is required by the protocol we want to use. Additional information about computing a checksum can be found in the section 3.1.2 on page 13.
- In the next step we send the data using the `write()` function as discussed earlier.
- Optionally we can `close()` the BPF device now, if we already sent all the data we want to transmit.

3.4 Winsock-API

As mentioned in section 2.1.1 on page 10 the Winsock-API does not allow any user interaction by itself, which would justify the usage of RAW-sockets. The rolled out Winsock-API is designed to be used with usual sockets using predefined protocols. Indeed, there is the possibility to use `SOCKET_RAW`, but it is mostly used for diagnostics. The limits, as also shown above are limiting to internet protocols IP, UDP and SCTP.

To write own protocols using Windows, or to receive packets, which are neither known or specified, there is the need to use libraries (npcap, libnet) to be able to support the Windows infrastructure. Otherwise the *invalid* or more *falsely marked as invalid* packets will be dropped by driver and OS.

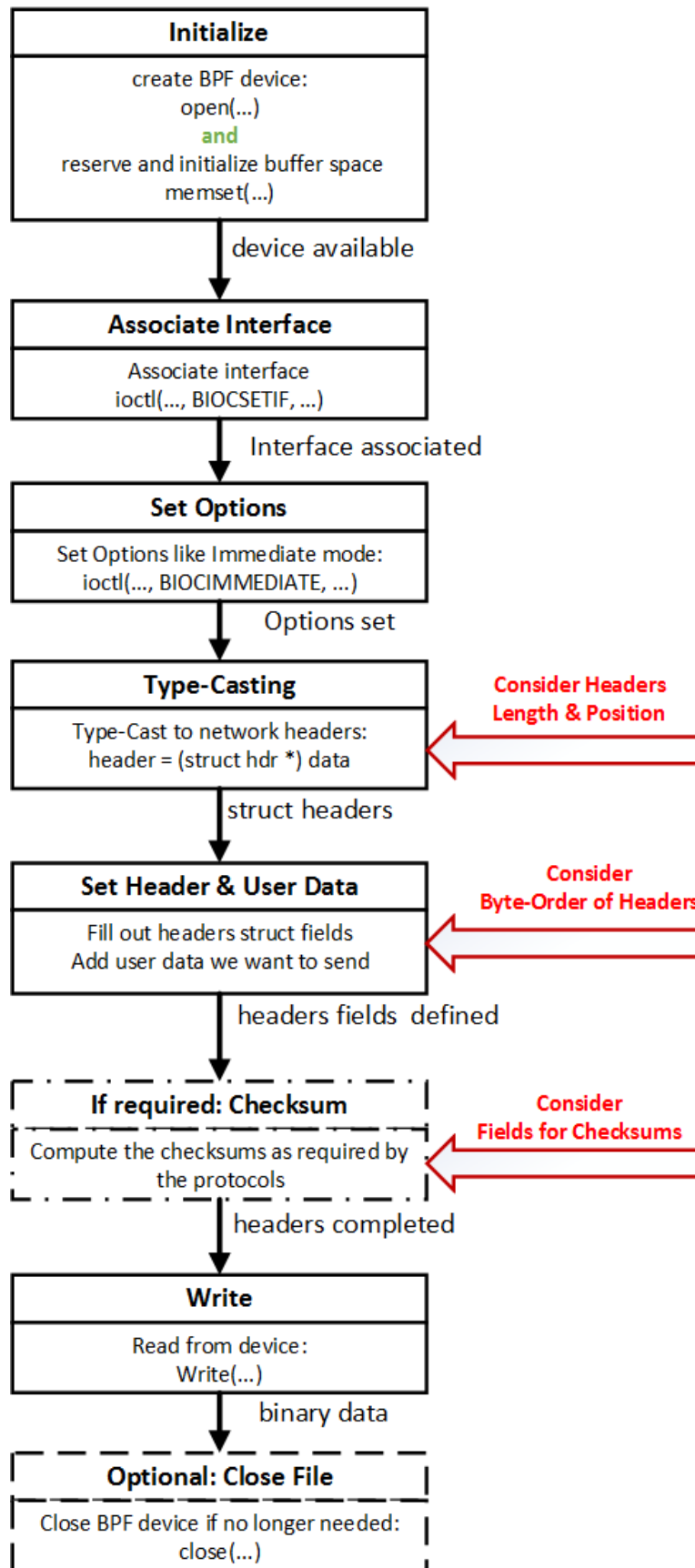


Figure 10: Structure of a BPF device layer 2 Write Operation

3.4.1 Preparations and Usage

If anyway there is no possibility to use the libraries, sockets of type `SOCK_RAW` can be created, despite being restricted as mentioned. A call to the socket function with *address family* set to `AF_INET` or `AF_IP6` and *type* to `SOCK_RAW` will return a raw ip socket. For the network layer one of the mentioned restricted protocols need to be set in *protocol* parameter. The following listing provides an extended example of this socket-function call. It also gives valid example values for the *type*, *family* and *protocol*. Complete description of the socket-function is found at msdn: [https://msdn.microsoft.com/en-us/library/windows/desktop/ms740506\(v=vs.85\).aspx](https://msdn.microsoft.com/en-us/library/windows/desktop/ms740506(v=vs.85).aspx)

See listing 13 in the appendix.

4 Programming with the libraries

This chapter will cover the two libraries that provide packet injection and capturing mechanisms for the most common operating systems. First, libnet, which is used for the packet injection, is introduced then we continue with the packet capture library, called pcap.

4.1 Libnet

Libnet is a library, mainly written in C, that defines an high-level portable interface for low-level network packet construction and injection. When libnet was designed the goal was to abstract out the pedantic architecture-specific details of low-level network packet tasks and provide some kind of platform-independent standard. It is under the BSD Licence and offers the ability to create and manipulate network packets from the link layer upwards. For each network layer libnet supports the manipulation of a set of protocols which are listed in figure 11. Libnet only supports the creation and sending of packets. It does not provide any functionalities to receive these packets. For latter purpose, one has to rely on an additional library called pcap.

The list of supported Operating Platforms is unfortunately no longer maintained. The latest state was that libnet supports the following OS [13]:

- Linux
- Windows
- FreeBSD
- OS X
- Solaris

The current version of libnet is 1.2. It is important to note that code from version 1.0.x will not work anymore as the syntax was simplified with version 1.1. Libnet was originally maintained on packetfactory.net by Mike D. Schiffman until 2004. However this page does not exist anymore and the author is also unreachable.

Since 2009 Sam Roberts maintains libnet on GitHub. As of today (December, 2016) the latest release is already four years ago but is still actively used [14].

4.1.1 Preparations

Libnet can be downloaded from the official libnet GitHub repository [14]. Alternatively, the currently latest version 1.2 is also available at sourceforge [15]. To install libnet, go to the directory, unzip it and run the following commands:

1. `./configure`
2. `make`
3. `make install` (might require root permissions)

Now `<libnet.h>` can be included in a C program. In order to successfully compile the program, run:

1. `gcc -Wall -g 'libnet-config --defines' -c foo.c`
2. `gcc -Wall foo.o -o output 'libnet-config --libs'`

4.1.2 Process of packet creation

In order to build and inject a network packet in libnet, there is a standard order of four operations:

1. **Library initialization**
2. **Packet Building**
3. **Packet Writing**
4. (Optional) **Packet destruction**

Library initialization The first step is to initialize the libnet library and to set up the environment. Doing so, the programmer receives a so called *libnet context*. This context maintains the state for the entire session which tracks all memory usage and packet construction. The *libnet context* is also often required as parameter for several functions such as packet building or injection.

Note that, the libnet library initialization can only be successfully executed with root permissions.

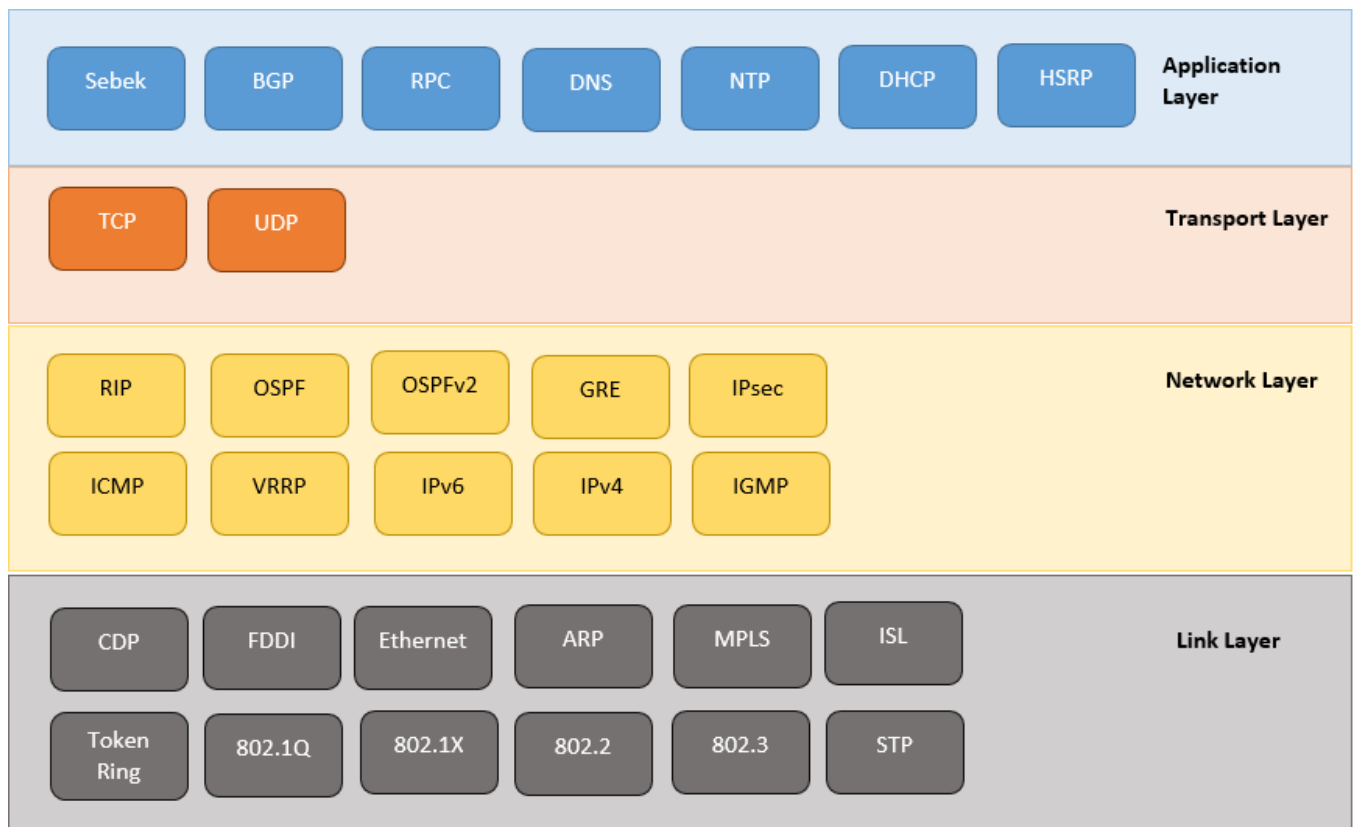


Figure 11: Overview of the packet construction in libnet

Packet building The second step is about the core functionality of libnet - the packet creation. In this phase the programmer calls a set of functions which are all of the structure `libnet_build_*`.

For every protocol header that libnet supports there is an own build function. Each `libnet_build()` function builds just a piece of the whole packet - namely the header of a single protocol. While it is possible to build an entire, ready-to-transmit packet with a single call to a `libnet_build()` function, normally more than one builder-class function call is required to construct a full packet. Every build function takes a series of arguments corresponding to the protocol header values as they appear on the wire. Thus, a programmer does not have to worry too much about low-level network details such as putting the bytes at the correct header position as this is done by the build function internally. Instead, he simply has to fill in the parameters of the function correctly. This process is very straight-forward but leads to `libnet_build()` functions which sometimes have a huge number of parameters which most of can be filled in with default values.

One important fact is that these build functions must be called in order, corresponding to how they should appear on the OSI layer model (from the highest protocol layer downwards). For example, to build a Network Time Protocol (NTP) packet, the programmer would call the `libnet_build()` functions in the following order:

1. `libnet_build_ntp()`
2. `libnet_build_udp()`
3. `libnet_build_ipv4()`
4. `libnet_build_ethernet()`

This process is also illustrated in figure 12. This ordering was not required until version 1.1.x where it became essential to properly link together the packet internally. Thus, source code of libnet before version 1.1.x is outdated and cannot be used anymore.

In case old source code shall be easily and quickly refactored one can refer to the `MIGRATION` file within the GitHub repository which gives instructions on how to update outdated source code.[16]

Packet write After having initialized the *libnet_session* and built the packet, the programmer can send it to the network.

Packet destruction The last step is optional. The programmer can decide to shut down the libnet session he created with the `init` function. Doing so, he closes the network interface and frees all internal memory structures used by the created packets.

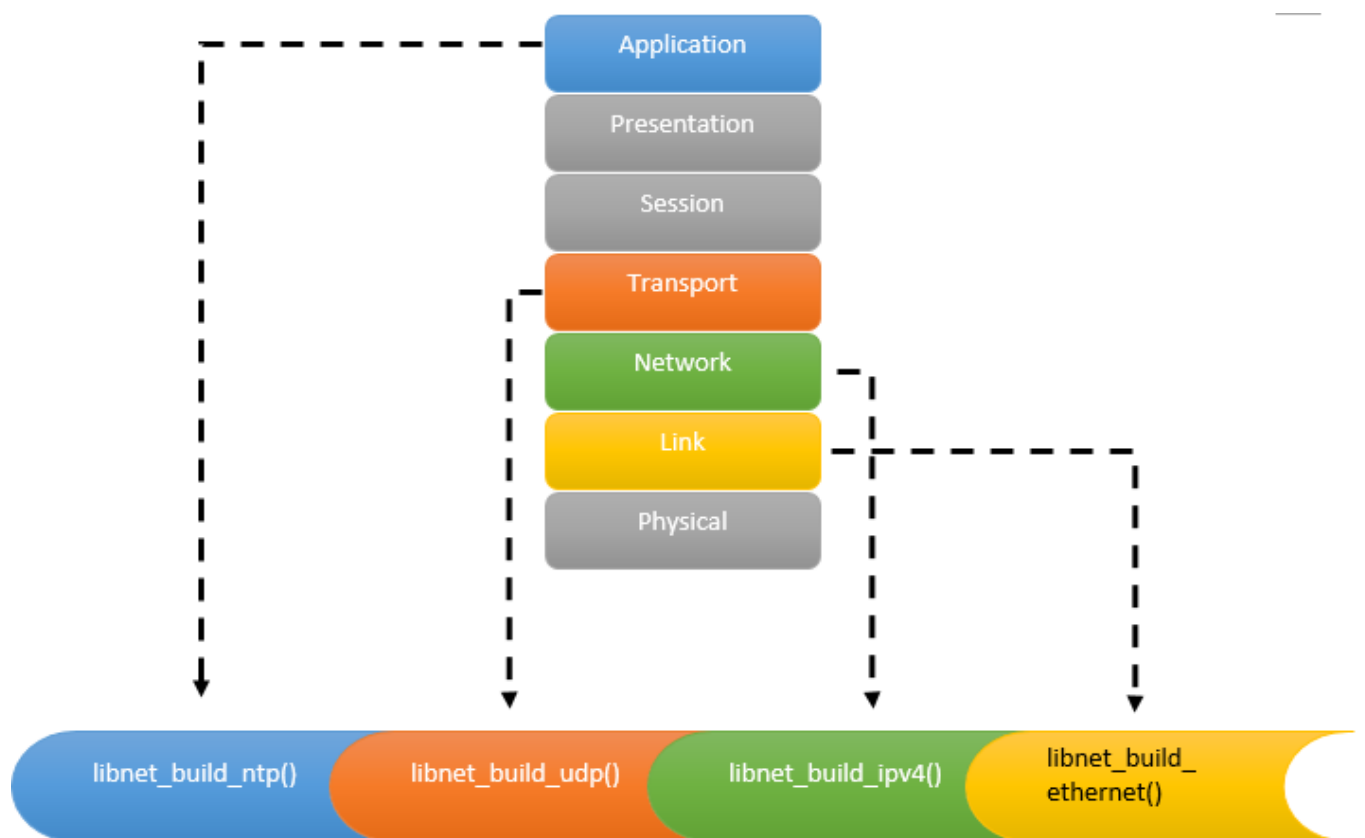


Figure 12: Overview of the packet construction in libnet

4.1.3 Functions

This section will have a deeper look at the exact syntax of the functions which perform the steps mentioned above. The definition of the functions is mainly retrieved from the official documentation of libnet.[17]

Library initialization The function to initialize the libnet environment is defined as follows:

```
libnet_t* libnet_init(int injection_type, char* device, char* err_buf)
```

Argument	Description
@SUCCESS	Returns libnet context.
@FAILURE	Returns NULL, err_buf will contain the reason.
int injection_type	The type of socket that will be used.
char* device	Use the specified network device for packet injection.
char* err_buf	Error message. Will only be filled in if function fails.

Table 21: Overview of `libnet_init()`

`libnet_init()` initializes the libnet library and creates the environment to work in. One important decision that the programmer has to make during this step is to decide whether libnet shall work with Link Layer sockets or Network-layer sockets (also called raw sockets). For this purpose the argument **injection_type** must be filled in with one of the following possible values: (values are all defined in libnet's header file[18])

- `LIBNET_LINK` for Link Layer sockets
- `LIBNET_RAW4` for raw sockets using IPv4
- `LIBNET_RAW6` for raw sockets using IPv6

If the programmer decides to use link layer sockets, he must later ensure to manually build a protocol header (such as Ethernet) in the link layer for his packet. In case he chooses raw sockets he can skip building a protocol header in the link layer and must only build headers down to the network layer.

In case of success `libnet_init()` returns the so-called libnet context which is important to save in a variable as it is required as an argument for nearly all the other functions of libnet.

For choosing the desired network device (*char** **device**) following formats are possible:

- Canonical string that references the device (such as **eth0** or **fxp0**).
- A dots and decimals representation of the device's IP address (such as 192.168.0.1).
- NULL. In this case the function will select the first device available.

For the creation of **err_buf** the constant `LIBNET_ERRBUF_SIZE` (defined in the libnet headerfile) can be used to ensure a fitting size of the character array.

Build Functions As already stated there is one build function for each protocol that libnet supports. Sometimes a protocol has different types which results in a protocol having more than one build function. For example for ICMP there are the following build functions available:

- `libnet_build_icmpv4_echo()`
- `libnet_build_icmpv4_mask()`
- `libnet_build_icmpv4_unreach()`
- `libnet_build_icmpv4_redirect()`
- `libnet_build_icmpv4_timeexceed()`
- `libnet_build_icmpv4_timestamp()`

As it would be simply too much to introduce each of the available build functions we will cover just three of the most frequently used. If the reader wants to learn more about how to use a specific build function he can refer to the **libnet-functions.h** file under `libnet/include/libnet/` where each build function is documented. For now, let's further investigate on how to build a UDP, IPv4 and an Ethernet packet.

The UDP build function is defined as follows:

```
libnet_ptag_t libnet_build_udp(u_int16_t sp, u_int16_t dp, u_int16_t len,
u_int16_t sum, u_int8_t* payload, u_int32_t payload_s, libnet_t* l, libnet_ptag_t
ptag)
```

Argument	Description
@SUCCESS	Returns a ptag identifier referring to the UDP packet.
@FAILURE	Returns -1 and libnet_get_error() returns the error reason.
<i>u_int16_t</i> sp	The UDP port of the source.
<i>u_int16_t</i> dp	The UDP port of the destination.
<i>u_int16_t</i> len	The length of the UDP packet (including payload).
<i>u_int16_t</i> sum	The checksum of the packet. If 0, libnet will autofill.
<i>u_int8_t*</i> payload	The optional payload. Can be NULL.
<i>u_int32_t</i> payload_s	The length of the payload. Is 0 if payload is NULL.
<i>libnet_t*</i> l	The pointer to a libnet context
<i>libnet_ptag_t</i> ptag	The protocol tag of the packet. It must be 0 if this is a new packet

Table 22: Overview of libnet_build_udp()

The IPv4 build function is defined as follows:

```
libnet_ptag_t libnet_build_ipv4(uint16_t ip_len, uint8_t tos, uint16_t id,
uint16_t frag, uint8_t ttl, uint8_t prot, uint16_t sum, uint32_t src, uint32_t dst,
const uint8_t* payload, uint32_t payload_s, libnet_t* l, libnet_ptag_t ptag)
```

Argument	Description
@SUCCESS	Returns a ptag identifier referring to the IPv4 packet.
@FAILURE	Returns -1 and libnet_get_error() can tell you why
<i>uint16_t</i> ip_len	The length of the IPv4 packet (including payload)
<i>uint8_t</i> tos	The type of service bits
<i>uint16_t</i> id	IP identification
<i>uint16_t</i> frag	Fragmentation bits
<i>uint8_t</i> ttl	Time to live
<i>uint8_t</i> prot	The upper layer protocol (above the ip header)
<i>uint16_t</i> sum	The checksum, 0 for libnet autofill
<i>uint32_t</i> src	The source IP address (in little endian)
<i>uint32_t</i> dst	The destination IP address (in little endian)
<i>const uint8_t*</i> payload	The optional payload. Can be NULL.
<i>uint32_t</i> payload_s	The length of the payload. Is 0 if payload is NULL.
<i>libnet_t*</i> l	The pointer to the libnet context.
<i>libnet_ptag_t</i> ptag	The protocol tag of the packet. It must be 0 if this is a new packet

Table 23: Overview of libnet_build_ipv4()

The Ethernet build function is defined as follows:

```
libnet_ptag_t libnet_build_ethernet(const uint8_t* dst, const uint8_t*
src, uint16_t type, const uint8_t* payload, uint32_t payload_s, libnet_t* l,
libnet_ptag_t ptag)
```

Argument	Description
@SUCCESS	Returns a ptag identifier referring to the IPv4 packet.
@FAILURE	-1 and libnet_get_error() can tell you why
<i>const uint8_t*</i> dst	The destination MAC address
<i>const uint8_t*</i> src	The source MAC address
<i>uint16_t</i> type	The upper layer protocol type
<i>const uint8_t*</i> payload	The optional payload. Can be NULL.
<i>uint32_t</i> payload_s	The length of the payload. Is 0 if payload is NULL.
<i>libnet_t*</i> l	The pointer to a libnet context.
<i>libnet_ptag_t</i> ptag	The protocol tag of the packet. It must be 0 if this is a new packet.

Table 24: Overview of libnet_build_ethernet()

When examining these three libnet_build() function examples, we discover the following:

All build functions return a protocol tag **ptag** on success and have an argument for this **ptag**. **Ptag** is used to identify a packet header after being created. The programmer must save this **ptag** if he intends to modify the packet later. To do so, he must call the same build function but with the respective **ptag** as argument. The libnet_build() function will then search for the existing packet and override the old values rather than creating a new one. In case the programmer wants to create a new packet he simply puts a "0" as argument for **ptag**.

Each build function contains an argument field for the libnet context which is usually the second last argument of the function. This field simply has to be filled in with the libnet context that is returned after the libnet_init() function.

Also, each build function has an argument for the payload that comes after the header data and, respectively, an argument for the size of that payload data. It is a simple interface to append arbitrary payloads to packets. It is completely up to the user to define which data exactly he or she wants to add to the packet header. Hence, the payload argument is a good possibility to implement an own type of protocol header. There is even an own build function for just adding raw payload data. It is defined as:

```
libnet_ptag_t libnet_build_data(const uint8_t* payload, uint32_t payload_s,  
libnet_t* l, libnet_ptag_t ptag)
```

The build_data() function can be used at any place to create additional, user-defined payload in between the other build functions.

Some functions contain an argument for the checksum. By default, checksums will be automatically calculated if "0" is put as parameter. In case it is desired to compute the checksum manually, please refer to function libnet_toggle_checksum() under section 4.1.3.

Some build functions require a length argument which indicates the whole length of the packet. The length of the whole packet is defined as the sum of the following components:

- Length of the header itself
- Length of the additional payload size
- Length of the encapsulated packets

For example, an IPv4 Packet which encapsulates a UDP header has the following size: the size of the IPv4 Header (20 bytes) + the size of its own payload (if it exists) + the size of the UDP header (8 bytes) + the size of the payload of UDP (if it exists).

In order to avoid to look up the size of each protocol header, libnet offers constants for the size of each header that it supports. They are saved in **libnet-headers.h** inside the libnet repository. The structure of the constant is LIBNET_name_of_protocol_H, e.g. LIBNET_UDP_H for the size of a UDP header or LIBNET_IPV4_H for the size of an IPv4 header respectively.

Some functions, such as libnet_build_ipv4() or libnet_build_ethernet(), require an argument that specifies the upper-layer protocol. Once again, there are constants available for this purpose. For example, IP Protocol types can be IPPROTO_ICMP, IPPROTO_UDP and IPPROTO_TCP. Libnet itself does not have a headerfile in its repository which defines these constants. Instead it retrieves them from a file called **netinet/in.h**. So, in order to retrieve a full list of available constants, please refer to this file. The same approach is used for the Ethernet protocol types which is required by the libnet_build_ethernet(). This time the respective headerfile with the constants is **netinet/if_ether.h**. If, for instance, the IP protocol is used above the link layer interface then the argument would be ETHERTYPE_IP.

There are also so-called auto-build functions. These are simply build function which require less parameters than the normal build functions. These functions are useful when one intends to build a header quickly because a granular level of control is not desired or required. Auto-build functions exist for the following header packets:

- Ethernet
- FDDI
- ARP
- IPv4
- IPv6
- Tokenring
- LinkLayer

As an example the autobuild function of IPv4 looks as follows:

```
libnet_ptag_t libnet_autobuild_ipv4(uint16_t len, uint8_t prot, uint32_t dst,
libnet_t* l)
```

Argument	Description
@SUCCESS	Returns a ptag identifier.
@FAILURE	Returns -1.
<i>uint16_t</i> len ,	Total length of the IP packet including all subsequent data
<i>uint8_t</i> prot	Upper layer protocol
<i>uint32_t</i> dst	Destination IPv4 address in little endian
<i>libnet_t*</i> l	Pointer to a libnet context

Table 25: Overview of libnet_autobuild_ipv4()

Comparing this function to the normal IPv4 we can see that the function takes the same length, protocol, and destination arguments as libnet_build_ipv4(). The function does not accept a **ptag** argument, but it does return a **ptag**. Thus, it can be used to build a new IP header but not to modify an existing one.

For more information on how the other autobuild functions look like, please refer to the **libnet-functions.h** file under the GitHub repository.[17]

Write The write function is defined as follows:

```
int libnet_write(libnet_t* l)
```

Argument	Description
@SUCCESS	Returns the number of bytes written
@FAILURE	Returns -1
<i>libnet_t*</i> l	Pointer to a libnet context

Table 26: Overview of libnet_write()

The function assumes that a *libnet_context* has been properly initialized with libnet_init() and that a previously constructed packet has been built inside this context via one or more calls to the libnet_build() functions. The write-function will then assemble the packet from the protocol blocks and send the packet either to the IP address for an injection at the LIBNET_RAW level or to a network hardware address if the LIBNET_LINK level was chosen during the initialization step.

Destruction The command to shut down a libnet context looks as follows:

```
void libnet_destroy(libnet_t* l)
```

It shuts down the libnet session referenced by **l**. It closes the network interface and frees all internal memory structures associated with **l**.

Other functions Besides the core functions introduced above, there are other interesting functions which can be useful. These are:

- Address resolution functions: When building the IP protocol header, it requires both a source and a destination IP address in little endian order with decimal values. As a programmer usually works with IP addresses in octet format and Big Endian order (such as 192.168.1.100) the addresses have to be converted beforehand. For this purpose, there are specific functions provided by libnet.

The definition of the function which converts from presentation format to address format is:

```
uint32_t libnet_name2addr4(libnet_t* l, char* host_name, uint8_t use_name)
```

Argument	Description
@SUCCESS	An IP number suitable for use with libnet_build() functions.
@FAILURE	Returns -1, which is technically “255.255.255.255”
libnet_t* l	The libnet context pointer
char* host_name	The presentation format address.
uint8_t use_name	Turn on/off address resolution (LIBNET_RESOLVE, LIBNET_DONT_RESOLVE)

Table 27: Overview of libnet_name2addr4()

This function takes a dotted decimal string (such as 192.168.19.1) or a canonical DNS name as argument for **host_name** and returns a decimal, little-endian ordered IPv4 address. In the official documentation, it says that the return value is network byte ordered. Network byte order however nearly always corresponds to big endian which is not the case here. A test of libnet_name2addr4() resulted in converting the IPv4 address 192.168.1.100 into 1677830336 which corresponds to a little endian order of previous IPv4 address. So the official documentation is misleading at this point.

If a DNS name is used, then LIBNET_RESOLVE must be used as argument for **use_name** so that a DNS lookup is performed. The function can fail if DNS lookup fails or if mode is set to LIBNET_DONT_RESOLVE and **host_name** refers to a canonical DNS name.

The counterpart to the previous function is the address to name conversion:

```
char* libnet_addr2name4(uint32_t in, uint8_t use_name)
```

Argument	Description
@SUCCESS	Returns a string of dots and decimals or a hostname
@FAILURE	This function cannot fail
uint32_t in	IPv4 address
uint8_t use_name	Turn on/off address resolution LIBNET_RESOLVE, LIBNET_DONT_RESOLVE

Table 28: Overview of libnet_addr2name4()

This function takes a little-endian ordered IPv4 address (e.g. the output of the previous function) and returns a string to either a canonical DNS name (if it has one) or a string of dotted decimals. This may incur a DNS lookup if the hostname and mode is set to LIBNET_RESOLVE. If mode is set to LIBNET_DONT_RESOLVE, no DNS lookup will be performed and the function will return a pointer to a dotted decimal string. The function cannot fail even if no canonical name exists. In this case it will simply return the dotted decimal string.

Address resolution functions do not only exist for IP addresses but also for MAC addresses. The function libnet_hex_aton() converts a colon separated hexadecimal address MAC address and returns a byte string suitable for use in a libnet_build() function.

```
uint8_t* libnet_hex_aton(const char* s, int* len)
```

The argument **s** corresponds to a MAC address such as 00:80:41:ae:fd:7e, and **len** defines the length of the MAC address.

- The possibility to let the checksum of a packet calculate automatically or manually. The respective function for this purpose is defined as follows:


```
int libnet_toggle_checksum(libnet_t* l, libnet_ptag_t ptag, int mode)
```

Argument	Description
@SUCCESS	Returns 1
@FAILURE	Returns -1
<i>libnet_t*</i> l	Pointer to a libnet context
<i>libnet_ptag_t</i> ptag	The ptag reference number
<i>int</i> mode	LIBNET_ON or LIBNET_OFF

Table 29: Overview of `libnet_toggle_checksum()`

By default, if a given protocol header contains the checksum field and it is set to "0", libnet will calculate the header checksum prior to injection. If the header is set to any other value, libnet will not calculate the header checksum. `libnet_toggle_checksum()` enables the programmer to over-ride this behavior. Mode `LIBNET_ON` switches auto-checksumming on for the specified ptag whereas `LIBNET_OFF` turns auto-checksumming off for the specified ptag. This assumes that the **ptag** of course refers to a protocol that has a checksum field. If the mode is set to `LIBNET_OFF`, libnet will clear the checksum flag and no checksum will be computed prior to injection. This assumes that the programmer will assign a value (zero or otherwise) to the checksum field. Often times this is useful if a precomputed checksum or some other predefined value is going to be used. Please note that when libnet is initialized with `LIBNET_RAW4`, the IPv4 header checksum will always be computed by the kernel prior to injection, regardless of what the programmer sets.

- Retrieve the total size of the constructed packet:

```
uint32_t libnet_getpacket_size(libnet_t* l)
```

This function returns the sum of the size of all of the headers that were built inside of **l**.

- Determine the name of the network device:

```
const char* libnet_getdevice(libnet_t* l)
```

This function will return the name of the network device which was used for packet injection. This function is especially useful if `NULL` was set as device parameter for the init function and the programmer wants to know which device has been selected now. Note it can be `NULL` without being an error. In this case there was no interface found.

It is not only possible to retrieve the name of network device. For instance, one can also retrieve the IP-address (both IPv4 and IPv6) of the device libnet was initialized with:

```
uint32_t libnet_get_ipaddr4(libnet_t* l)
```

The return value of `libnet_get_ipaddr4()` is in Big Endian order.

A similar function exists for the retrieval of the MAC address: `struct libnet_ether_addr*`
`libnet_get_hwaddr(libnet_t* l)`

- Retrieve information about statistics.

```
void libnet_stats(libnet_t* l, struct libnet_stats* ls)
```

Argument	Description
<i>libnet_t*</i> l	Pointer to a libnet context
<i>struct libnet_stats*</i> ls	Pointer to a libnet statistics structure

Table 30: Overview of `libnet_stats()`

In libnet there is a way to get statistics about the libnet context. These information cover the amount of packets written, amount of bytes written and the number of packet sending errors. These information will be retrieved in a structure which is passed as a pointer to the function. The structure is called `libnet_stats()` and is defined under *libnet-structures.h*.

- Retrieve the error message:

```
char* libnet_geterror(libnet_t* l)
```

Argument	Description
@SUCCESS	Returns an error string, NULL if none occurred
@FAILURE	This function cannot fail
<i>libnet_t*</i> l	Pointer to a libnet context

Table 31: Overview of `libnet_geterror()`

Returns the last error message that occurred during a `libnet_build()` function call. Remember there is a constant (`LIBNET_ERRBUF_SIZE`) that can be used to determine the optimal size for the array in which the error message shall be stored in.

4.2 libpcap

Libpcap is a portable library for network packet capturing written in C/C++. Similar to libnet, the goal of libpcap was to provide system-independent tools for packet filtering and capturing as almost every single OS-vendor had its own interface for this task. Captured packets via libpcap can be read from the data link layer upwards. Although it is written in C/C++ there are also wrappers for other common programming languages such as Java or Python.[19]

Libpcap is maintained by the developers behind tcpdump.org. It has an own GitHub repository where the whole source code is available as well as very detailed manpage where each function is documented.[20]

The latest release of libpcap is version 1.8.1 which was published in October 2016. Although libpcap considers itself to be system-independent it only supports UNIX-like operating systems including Linux, BSD and MacOS. For Windows there are special implementations of libpcap to ensure compatibility. Basically, they work very similar to libpcap however there are some minor differences. This topic will be covered in section 4.3. For now, we will focus on the basic libpcap implementation.

4.2.1 Preparation

To download libpcap go to the official site of tcpdump.org and unzip the file after the download.[21] To be able to use the functions provided by libpcap, include the headerfile **pcap/pcap.h** inside your C/C++ program. Libpcap requires root permission. Further details for each OS to take care of are mentioned in the README files in the GitHub repository.

4.2.2 Process of packet capturing

The structure of a libpcap program can be divided into five stages:

1. **Interface Selection**
2. **Initialize libpcap session**
3. **Define Filters**
4. **Initiate Sniffing process**
5. **Close libpcap session**

Interface Selection First of all the programmer must determine an interface where packets shall be captured. The selection can either be done manually by a user or automatically by libpcap. In latter case there is a function in libpcap which will simply choose the first network device that is suitable to be sniffed on. Interfaces in Linux usually have a name such as **eth0** or **wlan0**.

Initialize libpcap session In order to sniff for packets, a programmer must first initialize a libpcap session. During this initialization a network device is defined where the programmer actually wants to sniff for packets. However, it is also possible to listen to all available devices of the host at once. (on Linux systems with 2.2 or later kernels) Usually the initialized device corresponds with the selected interface in the previous step.

On top of that further decisions are made such as the maximum size that shall be received per packet or whether we want to capture all packets on the selected network interface or only those that are destined to the host. It is no problem to initialize more than one libpcap session for different network devices at the same time. The initialization step can be compared to the opening of a file for read/write operations.

Define Filters Sniffing on a network interface can result in receiving a lot of packets where many are not of interest for the user. Thus, there is the possibility to filter incoming packets. The filtering is done inside the kernel which makes the filtering process very efficient. Reason behind this that a kernel has to copy a received packet from kernel space to user space for further processing. This copy process however is very CPU intensive. With in-kernel filtering this process can be skipped. So if a programmer sets an in-kernel filter only the packets that pass the filter requirements will be copied to the user space resulting in a much faster performance. Setting filters is optional and not required if a programmer intends to receive all packets on the specified network device.

Every OS has its own packet filtering mechanism however many rely on the BSD Packet Filter (BPF) architecture. [22] Platforms that use this architecture are for example: BSD Operating Systems, MacOS and Linux. Libpcap also relies on BPF to filter its packets. In case a platform does not support BPF the packets can still be manually filtered by the programmer in the user space which results in much higher overhead.

Setting a filter in libpcap consists of three steps:

1. Constructing filter expression

2. Compiling filter expression into BPF program
3. Applying filter to libpcap session

Normally, one has to write a BPF program, whose language is similar to assembly, to define a filter expression. However libpcap provides a high-level language that abstracts the actual language and makes the creation of the filter much easier. The exact syntax of the BPF Filter expression will be covered in the functions section. For now, it is sufficient to know that the filter expression is just a simple string. As the computer cannot understand this, it must be compiled into a BPF program afterwards. In order for the filter to come into action, it must be applied to a libpcap session as the final step.

Initiate Sniffing process After having done all the pre-requirements one is now ready to start waiting for and receiving packets. By default, libpcap both captures the incoming as well as the outgoing packets of the host machine. There are two possibilities to capture packets in libpcap:

1. Possibility is to call a function that simply checks if there is a packet available and returns the content of this packet.
2. Possibility is to run a loop that processes each packet that will be received. It will not stop until a user defined amount of packets has been processed. It is also possible to let the loop run indefinitely.

Independently from what possibility the programmer chooses, the data inside the packets must be somehow read. However, when a packet is captured, all what the application has got, is just one big array of character bytes and along with that a few information such as the time of reception and the total size of captured bytes. For this reason the programmer must work through the single OSI layers starting from the Data Link Layer and identify the protocols used in each layer. It is important to note that the captured bytes always start from the data link layer, e.g. Ethernet or Wi-Fi 802.11.

In order to deal with the packets in a proper way the programmer must take care of the following:

- Identify the size of each header so that he or she can jump to the beginning of each header
- Cast the bytes to a struct which corresponds to the protocol header.
- In case the bytes of the packets shall be printed it must be taken care that the bytes are converted from network byte order to host byte order.

Identification of the header size Normally, the programmer expects to receive a specific type of packet because he has set a filter and therefore does not need to "blindly" read a packet. However, in some cases this might be inevitable. Then the following approach can be used to safely identify the single headers that are hiding inside the packet.

For the Data Link Layer the programmer can most of the time assume that it will be Ethernet but this is not always given since not all devices provide the same type of Data Link Layer headers. For this purpose libpcap provides a function (**pcap_datalink()**) that returns the Link Layer type of the device which was used for the initialization in the second step (Initialize libpcap session).

Assuming the Data Link layer type was Ethernet the identification of the upper layer protocols is straight-forward. The Ethernet header has a field called **ethertype** which is a 16-bit long value that specifies the upper-layer protocol such as IPv4. Table 32 contains a list of the three most common values. For a full list, please refer to:[23]

Network Layer Protocol	Ethertype Value
IPv4	0x0800
IPv6	0x86DD
Address Resolution Protocol (ARP)	0x0806

Table 32: Overview of Network Layer Protocols

If the network layer of the received packet is IPv4, the specification of the next higher layer (transport layer) is also easy. In this case there is a protocol field which contains a hexadecimal value. Again, the three most common ones can be found in table 33.[24]

Protocol	Value
ICMP	0x01
TCP	0x06
UDP	0x11

Table 33: Overview of Transport Layer Protocols

Casting Once a header has been identified, it is much more convenient to work with it when the bytes are cast so that they match with the structure of the header protocol. This way a programmer is able to easily access a specific field of the header. A programmer can decide to write its own structure in case he already knows how the header will look like or he can refer to predefined structures which usually exist for all common protocols, such as TCP, IPv4, or Ethernet. In Linux/Unix for instance there are a couple of header files already available that contain these common structures. For example, `netinet/tcp.h` has a structure for the TCP header. `netinet/if_ether.h` for Ethernet header.

Byte order The received bytes inside the packets are always in network byte order. When working with bytes, it is more convenient for them to be in host byte order. In the `arpa/inet.h` or `netinet/in.h` header files which are by default available in Linux/Unix, there are four functions which can be used for the conversion from network byte order into host byte order or vice versa. They are:

- `htonl()` and `htons()`: from host to network byte order whereas the function ending with *l* is for 32 bits and *s* for 16.
- `ntohs()` and `ntohl()`: from network to host byte order whereas the function ending with *l* is for 32 bits and *s* for 16.

During this step it might also be necessary to filter the packets in case one chooses to not use the BPF Filter for any reason such as the OS does not support BPF.

When working with raw packets one must always keep in mind that the packets might be malformed which means that the data inside is not reliable. Dealing with these malformed packets can quickly lead to errors such as segmentation fault. To avoid these errors as best as possible one can perform the following checks for a packet:

- Check the whole size of the packet and compare it to the expected size.
- If the packets are TCP/IP, one can check the checksum.
- Any data inside the packet that is intended to be used (like an IP address) should be checked beforehand.

Close libpcap session in the last step the libpcap session can be closed once sniffing is completed.

4.2.3 Functions

After having covered the basic flow of how a libpcap program works, we will now have a deeper look at the functions.

Interface Selection For the interface selection it might be the case that no function provided by libpcap is required if the user decides to set the device by himself. In general, it is better practice to let libpcap choose a suitable interface to ensure portability across different platforms. For this purpose, there is the function:

```
char* pcap_lookupdev(char* errbuf)
```

Argument	Description
@SUCCESS	Returns a network interface name
@FAILURE	Returns NULL, reason is in errbuf .
<i>char*</i> errbuf	Buffer in which in case of an error the reason is filled in

Table 34: Overview of `pcap_lookupdev()`

`pcap_lookupdev()` will return an interface suitable for packet sniffing which is not a loopback device. In case of an error the return value will be NULL and the error buffer will contain the reason. In the pcap headerfile (`pcap.h`) there is the constant `PCAP_ERRBUF_SIZE` which contains the optimal size for the error buffer.

As soon as the network device is selected, one can call another useful function which is:

```
int pcap_lookupnet(const char* device, bpf_u_int32* netp, bpf_u_int32* maskp,  
char* errbuf)
```

Argument	Description
@SUCCESS	Returns 0
@FAILURE	Returns -1, reason is in errbuf
<i>char*</i> device	Name of the network device
<i>bpf_u_int32*</i> netp	Buffer for the IP address
<i>bpf_u_int32*</i> maskp	Buffer for the network mask
<i>char*</i> errbuf	Buffer in which in case of an error the reason is filled in

Table 35: Overview of `pcap_lookupnet()`

This function saves the network IP address and the network mask of the selected device in two extra buffers. If it fails, it will fill in the error buffer just like the function `pcap_lookupdev()`. Especially the network mask buffer can be useful later when we want to set the filter so that it will only receive packets addressed to the host's network mask. Note that the order of the network mask and IP address is not in human readable form. See the code examples on how to convert it properly.

Another possibility is to retrieve a full list of all available devices with:

```
int pcap_findalldevs(pcap_if_t** alldevsp, char* errbuf)
```

Argument	Description
@SUCCESS	Returns 0
@FAILURE	Returns -1. Reason is inside errbuf
<i>pcap_if_t**</i> alldevsp	Pointer to a linked list containing all found devices
<i>char*</i> errbuf	Buffer to save the error message in case of failure

Table 36: Overview of `pcap_findalldevs()`

On success it returns 0 and it fills the buffer **alldevsp** with elements of type: *pcap_if_t* that represent one discovered device. It can occur that a program is not able to find all the available devices because it lacks of sufficient permission to open them for capturing. In this case they will not even be displayed in the list. The returned list of devices can also be NULL if no device has been found. The *pcap_if_t* structure contains four fields:

Member of Header	Description
<i>struct pcap_if*</i> next	Points to the next device in the list unless its NULL.
<i>char*</i> name	A character array containing the name of the device (e.g. eth0)
<i>char*</i> description	A human-readable description of the device. If there is none provided, it is NULL.
<i>struct pcap_addr*</i> addresses	A pointer to the first element of a list of network addresses for the device, it can be NULL if the device has no address.
<i>bpf_u_int32</i> flags	One of three possible device flags: PCAP_IF_LOOPBACK if device is a loopback interface, PCAP_IF_UP if device is up or PCAP_IF_RUNNING if the device is running.

Table 37: Overview of the *pcap_if_t* header

The field **address** itself has a specific structure type defined within **pcap.h**.^[25] It contains the following fields:

Member of Header	Description
<i>struct pcap_addr*</i> next	Points to the next device in the list unless its NULL.
<i>struct sockaddr*</i> addr	Points to a <i>struct sockaddr</i> containing an address.
<i>struct sockaddr*</i> netmask	Points to a <i>struct sockaddr</i> that contains the netmask corresponding to the address pointed to by addr
<i>struct sockaddr*</i> broadaddr	Points to a <i>struct sockaddr</i> that contains the broadcast address corresponding to the address pointed to by addr . It may be NULL if the device does not support broadcasts.
<i>struct sockaddr*</i> dstaddr	Points to a <i>struct sockaddr</i> that contains the destination address corresponding to the address pointed to by addr . It may be NULL if the device is not a point-to-point interface.

Table 38: Overview of the *pcap_addr* header

The structure *sockaddr* is defined in **netinet.h**/**in.h**:

Member of Header	Description
<i>unsigned short</i> sa_family	The address family (usually of type AF_XXX)
<i>char</i> sa_data[14]	14 bytes of protocol address

Table 39: Overview of the *sockaddr* header

It is important to note that the addresses in the list of addresses might be any type of address such as IPv4 or IPv6. For this reason it is necessary to check the *sa_family* member of the *struct sockaddr* before working with the content. IPv4 addresses have the *sa_family* value AF_INET, IPv6 addresses have AF_INET6. Also, one should not forget that the addresses are received in network byte order and might need to be converted to host byte order. If the data about the devices is no longer required it can be deleted with:

```
void pcap_freealldevs(pcap_if_t* alldevs)
```

where **alldevs** is a pointer to the list of devices that was returned by `pcap_findalldevs()`.

Initialize libpcap session To initialize a libpcap session the following function has to be called:

```
pcap_t* pcap_open_live(const char* device, int snaplen, int promisc, int to_ms,
char* errbuf)
```

Argument	Description
@SUCCESS	Returns <i>pcap_t*</i> handler
@FAILURE	Returns NULL
<i>const char*</i> device	Name of the network device
<i>int</i> snaplen	Specifies the snapshot length to be set on the handle
<i>int</i> promisc	Flag to set the promiscuous mode on/off
<i>int</i> to_ms	Read timeout in milliseconds
<i>char*</i> errbuf	Buffer in which in case of an error the reason is filled in

Table 40: Overview of `pcap_open_live()`

On success, the function returns a handler which is used for setting the filter and to start reading packets from the interface. Therefore it is important to save the return value of the function. If the function however fails, the error buffer will contain the reason and NULL is returned. The **device** argument contains the name of the selected device which was either returned by `pcap_lookupdev()` or which was entered manually. If the value NULL or **any** is chosen for this argument, then all packets from all interfaces are captured.

The **snaplen** argument specifies how many bytes shall be captured at most per packet. To ensure that the size is sufficient one can define a size of 2048 bytes. That should be sufficient for each captured packet. The promiscuous mode is turned off for the value 0 and turned on for each other value (usually 1). When the session is initialized in promiscuous mode it will also accept all packets that are not destined to the network card.

The argument **to_ms** specifies how many milliseconds the kernel should wait before copying captured information from kernel space to user space. 0 will wait until enough packets have arrived to the network interface. A list of common values are: 1000, 512, 10 or 0.

Define Filters As already outlined the definition and applying of filters consists of three steps. The first step is the construction of a filter expression which has a specific syntax.

One filter expression is made out of three parts called qualifiers. These qualifiers describe an identifier. An identifier is a value which usually is a IP address or a portnumber depending on the kind of qualifiers used. The three qualifiers that exist are:[26]

- **Type:** specifies the kind of thing the ID name or number refers to. Possible values are: **host**, **net**, **port** and **portrange**. Default is **host**.
- **Dir:** specifies a particular transfer direction to and/or from id. Possible values are: **src**, **dst**, **src or dst**, **src** and **dst**. Default is **src** or **dst** which means that all packets are captured that go to the identifier (**dst**) and all packets that come from the identifier (**src**). If IEEE 802.11 Wireless LAN link layers are used then the following headers are also valid: **ra**, **ta**, **addr1**, **addr2**, **addr3**, and **addr4**.
- **Proto:** this qualifier ensures that only packets of a specific protocol are captured. Possible protos are: **ether**, **fddi**, **tr**, **wlan**, **ip**, **ip6**, **arp**, **rarp**, **decnet**, **tcp** and **udp**. Default is that all protocols are captured.

The structure of the filter expression looks as follows: `[proto] + [dir] + [type] + identifier` whereas the identifier usually is an IP address or a port number. Note that not every qualifier must be used. If a qualifier is not specified in the filter expression the default values mentioned above are simply used.

To get a better understanding of how to construct a filter expression, here is a list of a few examples:

- `src host 192.168.1.77`: This filter returns packets whose source address is 192.168.1.71.
- `dst port 80`: This filter returns all packets that are sent to port 80.
- `ip[8]==5`: This filter returns packets whose IP TTL value equals 5.
- `tcp[13]==0x02 and (dst port 22 or dst port 23)`: This filter returns TCP packets with SYN Flag and whose destination is either port 22 or 23.

Note that the last example uses logical operators to combine multiple filter expressions. Allowed logical operator are: `and`, `or`, `not`. A filter expression is always saved in a character buffer. For more information about the construction of filter expression and more examples one can refer to the `pcap` filter man page.[26]

The next step is to compile that character buffer into the BPF format which `libpcap` is able to read. The particular function is:

```
int pcap_compile(pcap_t* p, struct bpf_program* fp, const char* str, int
optimize, bpf_u_int32 netmask)
```

Argument	Description
@SUCCESS	Returns 0
@FAILURE	Returns -1.
<i>pcap_t*</i> p	Pointer to <code>libpcap</code> session
<i>struct bpf_program*</i> fp	Buffer where the compiled program will be filled in
<i>const char*</i> str	Filter expression which shall be compiled
<i>int</i> optimize	Flag to optimize the resulting code or not.
<i>bpf_u_int32</i> netmask	IPv4 netmask on which packets shall be captured

Table 41: Overview of `pcap_compile()`

The first argument is the field for the pointer that is returned after having initiated a `libpcap` session successfully with `pcap_open_live`. The second argument is the buffer where the compiled code will be put in. The string argument is the filter expression. The `optimize` flag will optimize the compiled BPF program during compilation if set to true (any value other than 0). The last argument is the IPv4 netmask of the network on which the packets shall be captured. The `pcap_lookupnet()` function returns the netmask of a device in the appropriate type. If no network mask filtering is desired, then this parameter can simply set to `PCAP_NETMASK_UNKNOWN` which is a constant defined in the `pcap` header. It equals the netmask 255.255.255.255. On success the `fp` buffer will contain a compiled version of the filter. Otherwise the function `pcap_geterr()` can be called to learn about the failure.

With a compiled BPF filter we can apply the filter to the `pcap` session with:

```
int pcap_setfilter(pcap_t* p, struct bpf_program* fp)
```

Argument	Description
@SUCCESS	Returns 0
@FAILURE	Returns -1.
<i>pcap_t*</i> p	Pointer to <code>libpcap</code> session
<i>struct bpf_program*</i> fp	Buffer where the compiled program will be filled in

Table 42: Overview of `pcap_setfilter()`

The two arguments are the `fp` buffer which is the compiled version of the filter expression as well as the pointer to the `pcap` session to which the filter shall be applied. Once again, if anything fails, the error message can be found with the `pcap_geterr()` function.

Of course, it is possible to set more than one filter for one `pcap` session by calling the `pcap_setfilter` function more than once.

Initiate Sniffing process As already mentioned there are two possibilities to receive and process the incoming packets. The first one is to call a function which simply reads the next available packet and returns it. This function is:


```
const u_char* pcap_next(pcap_t* p, struct pcap_pkthdr* h)
```

Argument	Description
@SUCCESS	Returns the packet
@FAILURE	Returns NULL
<i>pcap_t*</i> p	Pointer to libpcap session
<i>struct pcap_pkthdr*</i> h	Pointer to packet header structure

Table 43: Overview of `pcap_next()`

The first argument is the libpcap session which the packet shall be read from. The second argument is the *pcap_pkthdr* structure which is defined within the pcap header file **pcap/pcap.h** as follows:

- *struct timeval* **ts**: the timestamp when the packet was received.
- *bpf_u_int32* **caplen**: the length of the information captured.
- *bpf_u_int32* **len**: the length of the total packet.

The difference between **caplen** and **len** is that **len** is the size of the full packet whereas **caplen** is only the part of the packet that was really captured and therefore is available. This could be the case if during the initialization of the pcap session a small value for the maximum size per packet was chosen.

On success the packet is returned as a string of character values. On failure NULL is returned. Unfortunately, there is no way to figure out the exact reason of failure. Reasons could be:

- Packet was discarded because it did not pass the filter
- The system has a standard read timeout which was triggered
- The system is non-blocking and there simply was no packet available at the time of the function call.

The second possibility to capture packets is with the help of the loop. This way, the capture does not need to be called over and over again if it is intended to receive a larger amount of packets. The respective function is *pcap_loop*:

```
int pcap_loop(pcap_t* p, int cnt, pcap_handler callback, u_char* user)
```

Argument	Description
@SUCCESS	Returns 0
@FAILURE	Returns -1
<i>pcap_t*</i> p	Pointer to libpcap session
<i>int</i> cnt	Number of packets to be captured
<i>pcap_handler</i> callback	Callback for each captured packet
<i>u_char*</i> user	Arguments for the callback function

Table 44: Overview of `pcap_loop()`

The first difference between the loop function and `pcap_next()` is that it does not directly return the packets' contents. Instead an integer value is returned which is 0 if the loop finished because the number of packets to be received represented by *cnt* was reached. Other return values are -1 if an error interrupted the loop or -2 if the loop was shut by a call of the function:

```
void pcap_breakloop(pcap_t* p).
```

(Argument is the handler of the pcap session) If `pcap_loop` shall run infinitely then **cnt** must be set to -1.

Since the return value does not provide the content of a packet anymore, there is another way how to retrieve the packet's content. This is the purpose of the callback argument which will be called every time a packet is ready to be read. The problem with callbacks is that the user cannot pass any arguments. This is the reason why the last argument (**user**) exists which enables the programmer to include arguments for the callback function. The callback function must have a specific prototype because otherwise `pcap_loop()` would not know how to deal with it. It looks as follows:

```
void function_name(u_char* userarg, const struct pcap_pkthdr* pkthdr, const u_char* packet)
```

Argument	Description
<i>u_char*</i> userarg	Pointer to the user arguments
<i>const struct pcap_pkthdr*</i> pkthdr	Pointer to packet header
<i>const u_char*</i> packet	Pointer to the captured packet

Table 45: Overview of callback

The first argument of the callback function is identical to the **user** pointer of `pcap_loop()`. The second one is the same packet header structure that is also used in `pcap_next()`. It contains information about the captured packet. The last argument is the captured packet itself. It is important to note that the **usearg** pointer is of type *u_char*. This means that the pointer must most likely be cast two times: One time when calling the `pcap_loop()` function and once again when using the argument inside the callback function. See the sample code section for an example which shows how to use `pcap_loop()` along with the user arguments. See line 143 of listing 16 in the appendix.

During the procession of the packets a programmer must additionally call the `pcap_datalink()` function in case he does not know for sure which type of packet was received. This function returns the link layer type of the packet. This way, the programmer can safely read the packet as he knows the structure of the received packet. The full specification of the function is: `int pcap_datalink(pcap_t* p)` where **p** is the handler of the pcap session. It returns an integer value which identifies the type of link layer header. Libpcap can distinguish more than 180 different link types. As it would be too much to list all of them here, the table 46 lists the two most common types. For a full list, please refer to URL.[27]

Data Link Type	Returned Integer value	Pcap Alias
Ethernet 10/100/1000 Mbs	1	DLT_EN10MB
Wi-Fi 802.11	6	DLT_IEEE802

Table 46: Common data link types

The field pcap alias refers to the name of constants representing the integer value which are defined within pcap.

Close Session Last but not least a programmer can close the session when no more packets shall be captured. This can be done with:

```
void pcap_close(pcap_t* p)
```

where **p** is the session handler. It frees all resources that were allocated by the session.

Other functions The functions above can be seen as libpcap's main routines that are mainly used. However there are some further functions which provide depending on the intended use useful functions. These include:

- Packet injection: Libpcap also provides packet injecting methods. However these are limited to two functions only where simply one large string is sent. In other words libpcap does not provide any functions to build the packets layer-by-layer as libnet does which makes packet injection in libpcap much harder. Therefore in our opinion it makes no sense to use the libpcap methods if you can also rely on libnet.
- Statistics: Libpcap provides a function which returns capture statistics of a pcap session. The function is defined as:

```
int pcap_stats(pcap_t* p, struct pcap_stat* ps)
```

The required structure *pcap_stat* is defined in `pcap/pcap.h`. It contains information like number of packets received/dropped or lost.

- Saving Packets: The content of a packet can be written to a so called *savefile* to preserve it for later use. The advantage of saving packets this way is that it can be read again with the same functions introduced above.

For a full list of all the available functions and their documentation, please refer to the official manpage of libpcap.[28]

A slightly modified example of an example-file, which shows the basic usage of the send-operation is in listing 14 in the appendix.

4.3 libpcap in Windows

When searching for libpcap under Windows one will quickly discover WinPcap. WinPcap is a special adaption of libpcap to also work on Windows. On its official site WinPcap even claims to be the "industry-standard windows

packet capture library".[29] The problem with Winpcap however is that its latest release was published in March 2013 (version 4.1.3). In this release the developers added support for Windows 8 and Windows Server 2012. However for the currently latest Windows 10 there is no official support which leads to the situation that WinPcap does not work on some builds under Windows 10 because the outdated NDIS5 API (which WinPcap relies on) has been removed. On top of that the developers of Winpcap stated in one of the last published announcements that they have little time to continue work on Winpcap. Although Winpcap considers itself as the industry standard, it is unsure whether the work on WinPcap will be continued in the near future.[29] For this reason, it makes more sense to focus on Npcap which is a codefork from WinPcap and still actively maintained. Npcap solves some of the disadvantages with Winpcap and provides new features including:

- Windows10 support: Npcap relies on the new NDIS 6 instead of the deprecated NDIS 5. Therefore Npcap also works on Windows 10.
- Raw 802.11 packet capture: Npcap is able to capture Wi-Fi packets which was not possible under WinPcap.
- Security: Npcap can be restricted so that only administrators can sniff packets.
- WinPcap Compatibility: Older programs that were written for WinPcap are still compatible and can be run in Npcap as well.
- Loopback Packet Capture: Npcap is able to capture packets from the loopback interface. For this purpose Npcap will create an adapter named Npcap Loopback Adapter after installation.
- Performance: Overall better performance than WinPcap.

4.3.1 Installation and Preparation

The installation of the npcap library is straightforward. An executable and an SDK can be found here: <https://nmap.org/npicap/>. The executable will place the required files on the system and a precompiled example (some can be found in the SDK examples folder) can be run. For developing applications, the SDK must be downloaded, extracted and be included into the C/C++ Project. Using an IDE, it is enough to point to the *.lib and *.h files for compiling and linking. Note: to be sure, that no packets are filtered by the Windows-Firewall it is recommended to disable it, since packets which are not usual TCP/IP packets, are seen as invalid and being dropped.

Installing the npcap executable is recommended here, because it will not interfere with a possible installation of *Wireshark*. Modified examples of an showing the basic usage of the send-operation is in listing 15 and showing the receive / dump option in listing 16 in the appendix.

4.3.2 Process of packet capturing

The process of packet capturing is entirely identical to the libnet process. Refer to 4.2.2.

4.3.3 Functions

The syntax of Npcap is identical to the syntax of Winpcap, e.g. Winpcap and Npcap provide the exact same functions. The syntax of WinPcap in turn relies on libpcap. So libpcap programs that are written under Linux/Unix are also portable to Windows. Because of this nearly all functions are identical to those that were introduced in the previous libpcap section. In fact, all of the above mentioned functions can be used in Winpcap as well. On top of that Npcap/WinPcap offer some extended functions that only work on Windows. These functions extend the capability of libpcap by providing remote packet capture, packet buffer size variation or high-precision packet injection. To learn about these functions in detail please visit:[4]

5 Literature

References

- [1] Linux man pages - raw-sockets. <http://man7.org/linux/man-pages/man7/raw.7.html>, 2014. visited on 2017-01-02.
- [2] Bpf manpage. [https://www.freebsd.org/cgi/man.cgi?bpf\(4\)](https://www.freebsd.org/cgi/man.cgi?bpf(4)), 2010. visited on 2017-01-21.
- [3] Tcp / ip raw sockets. [https://msdn.microsoft.com/en-us/library/windows/desktop/ms740548\(v=vs.85\).aspx](https://msdn.microsoft.com/en-us/library/windows/desktop/ms740548(v=vs.85).aspx). visited on 2017-03-15.
- [4] Winpcap's user manual. http://www.winpcap.org/docs/docs_412/html/group__wpcapfunc.html, 2009. visited on 2017-03-20.
- [5] André Volk. RAW Socket Programmierung und Einsatzfelder. Master's thesis, Universität Koblenz, the Netherlands, 2008.
- [6] Ip spoofing with bsd raw sockets interface. <http://www.enderunix.org/docs/en/rawipspoof/>, 2007. visited on 2017-01-21.
- [7] Using freebsd's bpf device with c/c++. <http://bastian.riek.ru/howtos/bpf/>, 2010. visited on 2017-01-21.
- [8] Linux man pages. <https://linux.die.net/man/>, 2014. visited on 2017-01-02.
- [9] Beej's network programming guide. http://beej.us/guide/bgnet/output/html/multipage/sockaddr_inman.html, 2017. visited on 2017-01-21.
- [10] Linux man pages - netdevice. <http://man7.org/linux/man-pages/man7/netdevice.7.html>, 2014. visited on 2017-01-02.
- [11] Micro howto. http://www.microhowto.info/howto/send_an_arbitrary_ipv4_datagram_using_a_raw_socket_in_c.html, 2017. visited on 2017-01-21.
- [12] Linux man pages - packet-sockets. <http://man7.org/linux/man-pages/man7/packet.7.html>, 2014. visited on 2017-01-02.
- [13] Libnet supported os. <https://github.com/sam-github/libnet/blob/master/libnet/doc/PORTED>, 2012. visited on 2017-03-12.
- [14] Libnet github repository. <https://github.com/sam-github/libnet>, 2016. visited on 2017-03-12.
- [15] Libnet download sourceforge. <https://sourceforge.net/projects/libnet-dev/>, 2014. visited on 2017-03-12.
- [16] Libnet migration instructions. <https://github.com/sam-github/libnet/blob/master/libnet/doc/MIGRATION>, 2013. visited on 2017-03-12.
- [17] Libnet function file. <https://github.com/sam-github/libnet/blob/master/libnet/include/libnet/libnet-functions.h>, 2013. visited on 2017-03-12.
- [18] Libnet header file. <https://github.com/sam-github/libnet/blob/master/libnet/include/libnet/libnet-headers.h>, 2012. visited on 2017-03-12.
- [19] Pcap supported libraries. <https://wiki.wireshark.org/Development/LibpcapFileFormat#Libraries>, 2015. visited on 2017-03-12.
- [20] Pcap github. <https://github.com/the-tcpdump-group/libpcap>, 2017. visited on 2017-03-12.
- [21] Tcpdump homepage. <http://www.tcpdump.org/>, 2017. visited on 2017-03-12.
- [22] S. McCanne and V. Jacobson. The bsd packet filter: A new architecture for user-level packet capture. <http://www.tcpdump.org/>, 1992. visited on 2017-03-12.
- [23] Ethernet types. <http://www.iana.org/assignments/ieee-802-numbers/ieee-802-numbers.xhtml#ieee-802-numbers-1>, 2017. visited on 2017-03-12.

- [24] Assigned internet protocol numbers. <https://www.iana.org/assignments/protocol-numbers/protocol-numbers.xhtml>, 2016. visited on 2017-03-12.
- [25] Pcap header file. <https://github.com/the-tcpdump-group/libpcap/blob/master/pcap/pcap.h>, 2017. visited on 2017-03-12.
- [26] Pcap filter mechanism. <http://www.tcpdump.org/manpages/pcap-filter.7.html>, 2015. visited on 2017-03-12.
- [27] Link-layer header types. <http://www.tcpdump.org/linktypes.html>, 2015. visited on 2017-03-12.
- [28] Pcap man page. <http://www.tcpdump.org/manpages/pcap.3pcap.html>, 2017. visited on 2017-03-12.
- [29] Winpcap homepage. <http://www.winpcap.org/>, 2013. visited on 2017-03-12.

A Appendix: Listings

A.1 rfc1071 checksum cpp

Listing 12: rfc1071 checksum cpp

```
1 unsigned short comp_chksum( unsigned short *addr, int len )
2 {
3     /**
4      * Quelle: RFC 1071
5      * Calculates the Internet-checksum
6      * Valid for the IP, ICMP, TCP or UDP header
7      *
8      * *addr : Pointer to the Beginning of the data
9      * (Checksummenfeld muss Null sein)
10     * len : length of the data (in bytes)
11     *
12     * Return : Checksum in network-byte-order
13     **/
14
15     long sum = 0;
16
17     while( len > 1 ) {
18         sum += *(addr++);
19         len -= 2;
20     }
21
22     if( len > 0 )
23         sum += * addr;
24
25     while (sum >> 16)
26         sum = ( ( sum & 0xffff ) + ( sum >> 16 ) );
27
28     sum = ~sum;
29
30     return ( ( u_short ) sum );
31 }
```

A.2 win socket cpp

Listing 13: win socket cpp

```
1  #ifndef UNICODE
2  #define UNICODE 1
3  #endif
4
5  // link with Ws2_32.lib
6  #pragma comment(lib, "Ws2_32.lib")
7
8  #include <winsock2.h>
9  #include <ws2tcpip.h>
10 #include <stdio.h>
11 #include <stdlib.h>    // Needed for _wtoi
12
13
14 int __cdecl wmain(int argc, wchar_t **argv)
15 {
16
17     //-----
18     // Declare and initialize variables
19     WSADATA wsaData = {0};
20     int iResult = 0;
21
22     //    int i = 1;
23
24     SOCKET sock = INVALID_SOCKET;
25     int iFamily = AF_UNSPEC;
26     int iType = 0;
27     int iProtocol = 0;
28
29     // Validate the parameters
30     if (argc != 4) {
31         wprintf(L"usage: %s <addressfamily> <type> <protocol>\n", argv[0]);
32         wprintf(L"socket opens a socket for the specified family, type, & protocol\n");
33         wprintf(L"%ws example usage\n", argv[0]);
34         wprintf(L"    %ws 0 2 17\n", argv[0]);
35         wprintf(L"    where AF_UNSPEC=0 SOCK_DGRAM=2 IPPROTO_UDP=17\n", argv[0]);
36         return 1;
37     }
38
39     iFamily = _wtoi(argv[1]);
40     iType = _wtoi(argv[2]);
41     iProtocol = _wtoi(argv[3]);
42
43     // Initialize Winsock
44     iResult = WSASStartup(MAKEWORD(2, 2), &wsaData);
45     if (iResult != 0) {
46         wprintf(L"WSAStartup failed: %d\n", iResult);
47         return 1;
48     }
49
50     wprintf(L"Calling socket with following parameters:\n");
51     wprintf(L"    Address Family = ");
52     switch (iFamily) {
53     case AF_UNSPEC:
54         wprintf(L"Unspecified");
55         break;
56     case AF_INET:
57         wprintf(L"AF_INET (IPv4)");
58         break;
59     case AF_INET6:
60         wprintf(L"AF_INET6 (IPv6)");
61         break;
62     case AF_NETBIOS:
63         wprintf(L"AF_NETBIOS (NetBIOS)");
64         break;
65     case AF_BTH:
66         wprintf(L"AF_BTH (Bluetooth)");
67         break;
68     default:
69         wprintf(L"Other");
70         break;
71     }
72     wprintf(L"    (%d)\n", iFamily);
```

```

73
74 wprintf(L" Socket type = ");
75 switch (iType) {
76 case 0:
77     wprintf(L"Unspecified");
78     break;
79 case SOCK_STREAM:
80     wprintf(L"SOCK_STREAM (stream)");
81     break;
82 case SOCK_DGRAM:
83     wprintf(L"SOCK_DGRAM (datagram)");
84     break;
85 case SOCK_RAW:
86     wprintf(L"SOCK_RAW (raw)");
87     break;
88 case SOCK_RDM:
89     wprintf(L"SOCK_RDM (reliable message datagram)");
90     break;
91 case SOCK_SEQPACKET:
92     wprintf(L"SOCK_SEQPACKET (pseudo-stream packet)");
93     break;
94 default:
95     wprintf(L"Other");
96     break;
97 }
98 wprintf(L" (%d)\n", iType);
99
100 wprintf(L" Protocol = %d = ", iProtocol);
101 switch (iProtocol) {
102 case 0:
103     wprintf(L"Unspecified");
104     break;
105 case IPPROTO_ICMP:
106     wprintf(L"IPPROTO_ICMP (ICMP)");
107     break;
108 case IPPROTO_IGMP:
109     wprintf(L"IPPROTO_IGMP (IGMP)");
110     break;
111 case IPPROTO_TCP:
112     wprintf(L"IPPROTO_TCP (TCP)");
113     break;
114 case IPPROTO_UDP:
115     wprintf(L"IPPROTO_UDP (UDP)");
116     break;
117 case IPPROTO_ICMPV6:
118     wprintf(L"IPPROTO_ICMPV6 (ICMP Version 6)");
119     break;
120 default:
121     wprintf(L"Other");
122     break;
123 }
124 wprintf(L" (%d)\n", iProtocol);
125
126 sock = socket(iFamily, iType, iProtocol);
127 if (sock == INVALID_SOCKET)
128     wprintf(L"socket function failed with error = %d\n", WSAGetLastError() );
129 else {
130     wprintf(L"socket function succeeded\n");
131
132     // Close the socket to release the resources associated
133     // Normally an application calls shutdown() before closesocket
134     // to disables sends or receives on a socket first
135     // This isn't needed in this simple sample
136     iResult = closesocket(sock);
137     if (iResult == SOCKET_ERROR) {
138         wprintf(L"closesocket failed with error = %d\n", WSAGetLastError() );
139         WSACleanup();
140         return 1;
141     }
142 }
143
144 WSACleanup();
145
146 return 0;
147 }

```


A.3 libnet tcp c

Listing 14: libnet tcp c

```
1 #include "libnet_test.h"
2
3 int
4 main(int argc, char *argv[]){
5     int c;
6     char *cp;
7     libnet_t *l;
8     libnet_ptag_t t;
9     char *payload;
10    u_short payload_s;
11    u_long src_ip, dst_ip;
12    u_short src_prt, dst_prt;
13    char errbuf[LIBNET_ERRBUF_SIZE];
14
15    printf("libnet 1.1 packet shaping: TCP + options[link]\n");
16
17    /*
18     * Initialize the library. Root privileges are required.
19     */
20    l = libnet_init(
21        LIBNET_LINK,                /* injection type */
22        NULL,                        /* network interface */
23        errbuf);                    /* error buffer */
24
25    if (l == NULL) {
26        fprintf(stderr, "libnet_init() failed: %s", errbuf);
27        exit(EXIT_FAILURE);
28    }
29
30    src_ip = 0;
31    dst_ip = 0;
32    src_prt = 0;
33    dst_prt = 0;
34    payload = NULL;
35    payload_s = 0;
36    while ((c = getopt(argc, argv, "d:s:p:")) != EOF){
37        switch (c)
38        {
39            /*
40             * We expect the input to be of the form 'ip.ip.ip.ip.port'. We
41             * point cp to the last dot of the IP address/port string and
42             * then separate them with a NULL byte. The optarg now points to
43             * just the IP address, and cp points to the port.
44             */
45            case 'd':
46                if (!(cp = strrchr(optarg, '.')))
47                {
48                    usage(argv[0]);
49                }
50                *cp++ = 0;
51                dst_prt = (u_short)atoi(cp);
52                if ((dst_ip = libnet_name2addr4(l, optarg, LIBNET_RESOLVE)) == -1)
53                {
54                    fprintf(stderr, "Bad destination IP address: %s\n", optarg);
55                    exit(EXIT_FAILURE);
56                }
57                break;
58            case 's':
59                if (!(cp = strrchr(optarg, '.')))
60                {
61                    usage(argv[0]);
62                }
63                *cp++ = 0;
64                src_prt = (u_short)atoi(cp);
65                if ((src_ip = libnet_name2addr4(l, optarg, LIBNET_RESOLVE)) == -1)
66                {
67                    fprintf(stderr, "Bad source IP address: %s\n", optarg);
68                    exit(EXIT_FAILURE);
69                }
70                break;
71            case 'p':
72                payload = optarg;
```

```

73         payload_s = strlen(payload);
74         break;
75     default:
76         exit(EXIT_FAILURE);
77     }
78 }
79
80 if (!src_ip || !src_prt || !dst_ip || !dst_prt)
81 {
82     usage(argv[0]);
83     exit(EXIT_FAILURE);
84 }
85
86 t = libnet_build_tcp_options(
87     (uint8_t*) "\003\003\012\001\002\004\001\011\010\012\077\077\077\077\000\000\000\000\000\000",
88     20,
89     1,
90     0);
91 if (t == -1)
92 {
93     fprintf(stderr, "Can't build TCP options: %s\n", libnet_geterror(1));
94     goto bad;
95 }
96
97 t = libnet_build_tcp(
98     src_prt, /* source port */
99     dst_prt, /* destination port */
100    0x01010101, /* sequence number */
101    0x02020202, /* acknowledgement num */
102    TH_SYN, /* control flags */
103    32767, /* window size */
104    0, /* checksum */
105    10, /* urgent pointer */
106    LIBNET_TCP_H + 20 + payload_s, /* TCP packet size */
107    (uint8_t*)payload, /* payload */
108    payload_s, /* payload size */
109    1, /* libnet handle */
110    0); /* libnet id */
111 if (t == -1)
112 {
113     fprintf(stderr, "Can't build TCP header: %s\n", libnet_geterror(1));
114     goto bad;
115 }
116
117 t = libnet_build_ipv4(
118    LIBNET_IPV4_H + LIBNET_TCP_H + 20 + payload_s, /* length */
119    0, /* TOS */
120    242, /* IP ID */
121    0, /* IP Frag */
122    64, /* TTL */
123    IPPROTO_TCP, /* protocol */
124    0, /* checksum */
125    src_ip, /* source IP */
126    dst_ip, /* destination IP */
127    NULL, /* payload */
128    0, /* payload size */
129    1, /* libnet handle */
130    0); /* libnet id */
131 if (t == -1)
132 {
133     fprintf(stderr, "Can't build IP header: %s\n", libnet_geterror(1));
134     goto bad;
135 }
136
137 t = libnet_build_ethernet(
138    enet_dst, /* ethernet destination */
139    enet_src, /* ethernet source */
140    ETHERTYPE_IP, /* protocol type */
141    NULL, /* payload */
142    0, /* payload size */
143    1, /* libnet handle */
144    0); /* libnet id */
145 if (t == -1)
146 {
147     fprintf(stderr, "Can't build ethernet header: %s\n", libnet_geterror(1));
148     goto bad;

```

```

149     }
150
151     /*
152      * Write it to the wire.
153      */
154     c = libnet_write(l);
155     if (c == -1)
156     {
157         fprintf(stderr, "Write error: %s\n", libnet_geterror(l));
158         goto bad;
159     }
160     else
161     {
162         fprintf(stderr, "Wrote %d byte TCP packet; check the wire.\n", c);
163     }
164
165     libnet_destroy(l);
166     return (EXIT_SUCCESS);
167 bad:
168     libnet_destroy(l);
169     return (EXIT_FAILURE);
170 }
171
172 void
173 usage(char *name)
174 {
175     fprintf(stderr,
176         "usage: %s -s source_ip.source_port -d destination_ip.destination_port"
177         " [-p payload]\n",
178         name);
179 }

```

A.4 sendpack c

Listing 15: sendpack c

```
1 #include <stdlib.h>
2 #include <stdio.h>
3 #include <pcap.h>
4
5 int main(int argc, char **argv){
6     pcap_t *fp;
7     char errbuf[PCAP_ERRBUF_SIZE];
8     u_char packet[100];
9     int i;
10
11     /* Check the validity of the command line */
12     if (argc != 2){
13         printf("usage: %s interface", argv[0]);
14         return 1;
15     }
16
17     /* Open the adapter */
18     if ((fp = pcap_open_live(
19         argv[1], // name of the device
20         65536, // portion of the packet to capture. It doesn't matter in this case
21         1, // promiscuous mode (nonzero means promiscuous)
22         1000, // read timeout
23         errbuf // error buffer
24     )) == NULL){
25         fprintf(stderr, "\nUnable to open the adapter. %s is not supported by WinPcap\n", argv[1]);
26         return 2;
27     }
28
29     /* Setting DMAC to B4 - B6 - 76 - D1 - 0F - C9
30      This is an example and does not need to be a valid MAC */
31     packet[0]=0xB4;
32     packet[1]=0xB6;
33     packet[2]=0x76;
34     packet[3]=0xD1;
35     packet[4]=0x0F;
36     packet[5]=0xC9;
37
38     /* set SMAC to 2:2:2:2:2:2 */
39     packet[6]=2;
40     packet[7]=2;
41     packet[8]=2;
42     packet[9]=2;
43     packet[10]=2;
44     packet[11]=2;
45
46     /* Fill the rest of the packet with data */
47     for(i = 12; i < 100; i++){
48         packet[i]= (u_char)i;
49     }
50
51     /* Send down the packet */
52     if (pcap_sendpacket(
53         fp, // Adapter
54         packet, // buffer with the packet
55         100 // size
56     ) != 0){
57         fprintf(stderr, "\nError sending the packet: %s\n", pcap_geterr(fp));
58         return 3;
59     }
60     printf("Packet sent.");
61     pcap_close(fp);
62     return 0;
63 }
```

A.5 dump c

Listing 16: dump c

```
1  #ifndef _MSC_VER
2  /*
3   * we do not want the warnings about the old deprecated and unsecure CRT functions
4   * since these examples can be compiled under *nix as well
5   */
6  #define _CRT_SECURE_NO_WARNINGS
7  #endif
8
9  #include "pcap.h"
10
11 /* 4 bytes IP address */
12 typedef struct ip_address{
13     u_char byte1;
14     u_char byte2;
15     u_char byte3;
16     u_char byte4;
17 }ip_address;
18
19 /* IPv4 header */
20 typedef struct ip_header{
21     u_char ver_ihl; // Version (4 bits) + Internet header length (4 bits)
22     u_char tos; // Type of service
23     u_short tlen; // Total length
24     u_short identification; // Identification
25     u_short flags_fo; // Flags (3 bits) + Fragment offset (13 bits)
26     u_char ttl; // Time to live
27     u_char proto; // Protocol
28     u_short crc; // Header checksum
29     ip_address saddr; // Source address
30     ip_address daddr; // Destination address
31     u_int op_pad; // Option + Padding
32 }ip_header;
33
34 /* UDP header*/
35 typedef struct udp_header{
36     u_short sport; // Source port
37     u_short dport; // Destination port
38     u_short len; // Datagram length
39     u_short crc; // Checksum
40 }udp_header;
41
42 /* UDP header*/
43 typedef struct eth_header{
44     u_char hwdes[6]; // Destination MAC
45     u_char hwsrc[6]; // Source MAC
46     u_char etherType[2]; // Ether Type 08 00 for ipv4
47 }eth_header;
48
49
50 enum FORMAT {BIN, DEC, HEX}; //For custom print function
51
52 /* prototype of the packet handler */
53 void packet_handler(u_char *param, const struct pcap_pkthdr *header, const u_char *pkt_data);
54
55 /* protortypes for printfunctions */
56 void print_bytes(char* text, void *object, size_t size, enum FORMAT f);
57 void print_bytes_noInc(char* text, void *object, size_t size, enum FORMAT f);
58
59 int main(){
60     struct sockaddr_in *saServer;
61     char errbuf[PCAP_ERRBUF_SIZE];
62     char packet_filter[] = ""; //kernel packet filter, "ip and udp" would be valid, iff only valid ip/udp
        packets are sent
63     u_int netmask;
64     pcap_t *adhandle; //for initializing device, whill hold the chosen one
65     pcap_if_t *alldevs; //for initializing device, will hold all
66     pcap_if_t *d; //for initializing device, iterator
67     int inum, i=0; //for initializing device, iterator for printing, choosing
68     struct bpf_program fcode;
69
70
71     /* Retrieve the device list */
```

```

72 if(pcap_findalldevs(&alldevs, errbuf) == -1){
73     fprintf(stderr, "Error in pcap_findalldevs: %s\n", errbuf);
74     exit(1);
75 }
76
77 /* Print the list */
78 for(d=alldevs; d; d=d->next){
79     printf("%d. %s", ++i, d->name);
80     if (d->description)
81         printf(" (%s)\n", d->description);
82     else
83         printf(" (No description available)\n");
84 }
85
86 if(i==0){
87     printf("\nNo interfaces found! Make sure WinPcap is installed.\n");
88     return -1;
89 }
90
91 printf("Enter the interface number (1-%d):", i);
92 scanf("%d", &inum);
93
94 /* Check if the user specified a valid adapter */
95 if(inum < 1 || inum > i){
96     printf("\nAdapter number out of range.\n");
97     pcap_freealldevs(alldevs); /* Free the device list */
98     return -1;
99 }
100
101 /* Jump to the selected adapter */
102 for(d=alldevs, i=0; i< inum-1 ;d=d->next, i++);
103
104 /* Open the adapter */
105 if ((adhandle= pcap_open_live(d->name, // name of the device
106     65536, // portion of the packet to capture: 65536 grants that the whole packet will be captured on
107         all the MACs.
108     1, // promiscuous mode (nonzero means promiscuous)
109     1000, // read timeout
110     errbuf // error buffer
111     )) == NULL){
112     fprintf(stderr, "\nUnable to open the adapter. %s is not supported by WinPcap\n");
113     pcap_freealldevs(alldevs);
114     return -1;
115 }
116
117 if(d->addresses != NULL)
118     /* Retrieve the mask of the first address of the interface */
119     netmask=((struct sockaddr_in *) (d->addresses->netmask))->sin_addr.S_un.S_addr;
120 else
121     /* If the interface is without addresses we suppose to be in a C class network */
122     netmask=0xffffffff;
123
124 /* compile the filter */
125 if (pcap_compile(adhandle, &fcode, packet_filter, 1, netmask) < 0 ){
126     fprintf(stderr, "\nUnable to compile the packet filter. Check the syntax.\n");
127     pcap_freealldevs(alldevs);
128     return -1;
129 }
130
131 /* set the filter */
132 if (pcap_setfilter(adhandle, &fcode)<0){
133     fprintf(stderr, "\nError setting the filter.\n");
134     pcap_freealldevs(alldevs);
135     return -1;
136 }
137
138 printf("\nlistening on %s...\n", d->description);
139
140 /* At this point, we don't need any more the device list. Free it */
141 pcap_freealldevs(alldevs);
142
143 /* start the capture */
144 pcap_loop(adhandle, 0, packet_handler, NULL);
145 return 0;
146 }

```

```

147 /* Callback function invoked by libpcap for every incoming packet */
148 void packet_handler(u_char *param, const struct pcap_pkthdr* header, const u_char *pkt_data){
149
150     /* buffers for time */
151     struct tm *ltime;
152     char timestr[16];
153     time_t local_tv_sec;
154
155     /* buffers for header data */
156     ip_header *ih;
157     udp_header *uh;
158     eth_header* eh;
159     u_int ip_len;
160     u_short sport, dport;
161
162     (VOID)(param); //unused parameter
163
164     /* convert the timestamp to readable format */
165     local_tv_sec = header->ts.tv_sec;
166     ltime = localtime(&local_tv_sec);
167     strftime( timestr, sizeof timestr, "%H:%M:%S", ltime);
168
169     /* print timestamp and length of the packet */
170     printf("%s.%.6d len:%d Ethernet Data:\n", timestr, header->ts.tv_usec, header->len);
171     print_bytes("Whole Data", pkt_data, header->len, HEX);
172     eh = (eth_header*)pkt_data;
173
174     /* Ethernet (DLL) Header */
175     print_bytes("\nMAC Dest\t\t", &(eh->hwdes), 6, HEX);
176     print_bytes("MAC Source\t\t", &(eh->hwsrc), 6, HEX);
177     print_bytes("Ethertype\t\t", &(eh->etherType), 2, HEX);
178
179     ih = (ip_header *) (pkt_data + 14); //length of ethernet header
180
181     /* position of udp (Network) header */
182     ip_len = (ih->ver_ihl & 0xf) * 4;
183     uh = (udp_header *) ((u_char*)ih + ip_len);
184
185     /* convert network to host byte order */
186     sport = ntohs( uh->sport );
187     dport = ntohs( uh->dport );
188
189
190     /* Transport (IP) Header */
191     print_bytes_noInc("\nVersion (H), IHL (L)\t", &(ih->ver_ihl), 1, BIN);
192     print_bytes_noInc("TypeOfService DSCP:6; ECN:2", &(ih->tos), 1, BIN);
193     print_bytes_noInc("Total Lenth\t\t", &(ih->tlen), 2, DEC);
194     print_bytes_noInc("Identification\t\t", &(ih->identification), 2, HEX);
195     print_bytes_noInc("Flags:3, FragOff:13\t", &(ih->flags_fo), 2, BIN);
196     print_bytes_noInc("TTL\t\t\t", &(ih->tttl), 1, DEC);
197     print_bytes_noInc("Protocol\t\t", &(ih->proto), 1, BIN);
198     print_bytes_noInc("Checksum\t\t", &(ih->crc), 2, DEC);
199
200     /* print ip addresses and udp ports */
201     printf("Source IP \t\t\t\t[ %d.%d.%d.%d ]\n",
202         ih->saddr.byte1,
203         ih->saddr.byte2,
204         ih->saddr.byte3,
205         ih->saddr.byte4);
206     printf("Destination IP\t\t\t\t[ %d.%d.%d.%d ]\n",
207         ih->daddr.byte1,
208         ih->daddr.byte2,
209         ih->daddr.byte3,
210         ih->daddr.byte4);
211     print_bytes_noInc("\nSource Port\t\t", &sport, 2, DEC);
212     print_bytes_noInc("Destination Port\t", &dport, 2, DEC);
213     print_bytes_noInc("datagram Length\t\t", &(uh->len), 2, DEC);
214
215     printf("\n\n");
216 }
217
218 void print_bytes(char* text, void *object, size_t size, enum FORMAT f){
219     /* buffers */
220     const u_char * const bytes = object;
221     size_t i;
222

```

```

223 printf("%s \t[ ", text);
224 switch (f){
225     case BIN:
226         for (i = 0; i < size; i++){
227             for (int j = 128; j > 0; j >= 1) {
228                 printf("%d", (bytes[i] & j) == j ? 1 : 0);
229                 if(j == 16)
230                     printf(" ");
231             }
232             printf(" ");
233         }
234         break;
235     case DEC:
236         if (size == 1){
237             printf("%d ", bytes[0]);
238             break;
239         }
240         for (i = 0; i+1 < size; i+=2)
241             printf("%d ", bytes[i + 1] + bytes[i] * 265);
242         break;
243     case HEX:
244         for (i = 0; i < size; i++){
245             printf("%02X ", bytes[i]);
246         }
247         break;
248     default:
249         break;
250 }
251 printf("]\n");
252 }
253
254 /* this function will not increase the pointer */
255 void print_bytes_noInc(char* text, void *object, size_t size, enum FORMAT f) {
256     void* obj = object; //copying the data, and passing
257     print_bytes(text, obj, size, f);
258 }

```


B Appendix: Tables

B.1 Protocol Types of <netinet/in.h>

Constant	Description
IPPROTO_IP	Dummy protocol.
IPPROTO_HOPOPTS	IPv6 Hop-by-Hop options.
IPPROTO_ICMP	Internet Control Message Protocol.
IPPROTO_IGMP	Internet Group Management Protocol.
IPPROTO_IPIP	IPIP tunnels (older KA9Q tunnels use 94).
IPPROTO_TCP	Transmission Control Protocol.
IPPROTO_EGP	Exterior Gateway Protocol.
IPPROTO_PUP	PUP protocol.
IPPROTO_UDP	User Datagram Protocol.
IPPROTO_IDP	XNS IDP protocol.
IPPROTO_TP	SO Transport Protocol Class 4.
IPPROTO_IPV6	IPv6 header.
IPPROTO_ROUTING	IPv6 routing header.
IPPROTO_FRAGMENT	IPv6 fragmentation header.
IPPROTO_RSVP	Reservation Protocol.
IPPROTO_GRE	General Routing Encapsulation.
IPPROTO_ESP	Encapsulating security payload.
IPPROTO_AH	Authentication header.
IPPROTO_ICMPV6	ICMPv6.
IPPROTO_NONE	IPv6 no next header.
IPPROTO_DSTOPTS	IPv6 destination options.
IPPROTO_MTP	Multicast Transport Protocol.
IPPROTO_ENCAP	Encapsulation Header.
IPPROTO_PIM	Protocol Independent Multicast.
IPPROTO_COMP	Compression Header Protocol.
IPPROTO_SCTP	Stream Control Transmission Protocol.
IPPROTO_RAW	Raw IP packets.
IPPROTO_MAX	No description.

Table 47: Protocol Types defined in <netinet/in.h> [8]

B.2 Linux Protocol Types defined in linux if_ether.h

Flag	Description
ETH_P_LOOP	Ethernet Loopback packet
ETH_P_PUP	Xerox PUP packet
ETH_P_PUPAT	Xerox PUP Addr Trans packet
ETH_P_IP	Internet Protocol packet
ETH_P_X25	CCITT X.25
ETH_P_ARP	Address Resolution packet
ETH_P_BPQ	G8BPQ AX.25 Ethernet Packet [not officially registered]
ETH_P_IEEE8023PUP	Xerox IEEE802.3 PUP packet
ETH_P_IEEE8023PUPAT	Xerox IEEE802.3 PUP Address Transport packet
ETH_P_DEC	DEC Assigned protocol
ETH_P_DNA_DL	DEC DNA Dump/Load
ETH_P_DNA_RC	DEC DNA Remote Console
ETH_P_DNA_RT	DEC DNA Routing
ETH_P_LAT	DEC LAT
ETH_P_DIAG	DEC Diagnostics
ETH_P_CUST	DEC Customer use
ETH_P_SCA	DEC Systems Communications Architecture
ETH_P_RARP	Reverse Address Resolution packet
ETH_P_ATALK	Appletalk DDP
ETH_P_AARP	Appletalk AARP
ETH_P_8021Q	802.1Q VLAN Extended Header
ETH_P_IPX	IPX over DIX
ETH_P_IPV6	IPv6 over bluebook
ETH_P_WCCP	Web-cache coordination protocol
ETH_P_PPP_DISC	PPPoE discovery messages
ETH_P_PPP_SES	PPPoE session messages
ETH_P_MPLS_UC	MPLS Unicast traffic
ETH_P_MPLS_MC	MPLS Multicast traffic
ETH_P_ATMMP	MultiProtocol Over ATM
ETH_P_ATMFATE	Frame-based ATM Transport over Ethernet
ETH_P_AOE	ATA over Ethernet
ETH_P_802_3	Dummy type for 802.3 frames
ETH_P_AX25	Dummy protocol id for AX.25
ETH_P_ALL	Every packet
ETH_P_802_2	802.2 frames
ETH_P_SNAP	Internal only
ETH_P_DDCMP	DEC DDCMP: Internal only
ETH_P_WAN_PPP	Dummy type for WAN PPP frames
ETH_P_PPP_MP	Dummy type for PPP MP frames
ETH_P_LOCALTALK	Localtalk pseudo type
ETH_P_PPPTALK	Dummy type for Atalk over PPP
ETH_P_TR_802_2	802.2 frames
ETH_P_MOBITEX	Mobitex
ETH_P_CONTROL	Card specific control frames
ETH_P_IRDA	Linux-IrDA
ETH_P_ECONET	Acorn Econet
ETH_P_HDLC	HDLC frames
ETH_P_ARCNET	ArcNet

Table 48: Linux Protocol Types defined in
<linux/if_ether.h>

B.3 Socket level options for `setsockopt()`

Flag	Description
SO_ACCEPTCONN	Indicates whether or not this socket has been marked to accept connections.
SO_BINDTODEVICE	Bind this socket to a particular device like 'eth0'.
SO_BROADCAST	Set or get the broadcast flag.
SO_BSDCOMPAT	Enable BSD bug-to-bug compatibility. If enabled ICMP errors received for a UDP socket will not be passed to the user program.
SO_DEBUG	Enable socket debugging. Only allowed for processes with the CAP_NET_ADMIN capability or an effective user ID of 0.
SO_DONTROUTE	Don't send via a gateway, only send to directly connected hosts. The same effect can be achieved by setting the MSG_DONTROUTE flag on a socket send() operation. Expects an integer boolean flag.
SO_KEEPAIVE	Enable sending of keep-alive messages on connection-oriented sockets. Expects an integer Boolean flag.
SO_LINGER	When enabled, a close() or shutdown() will not return until all queued messages for the socket have been successfully sent or the linger timeout has been reached.
SO_MARK	Set the mark for each packet sent through this socket.
SO_OOBINLINE	If this option is enabled, out-of-band data is directly placed into the receive data stream.
SO_PASSCRED	Enable or disable the receiving of the SCM_CREDENTIALS control message
SO_PRIORITY	Set the protocol-defined priority for all packets to be sent on this socket. Linux uses this value to order the networking queues. Setting a priority outside the range 0 to 6 requires the CAP_NET_ADMIN capability.
SO_RCVBUF	Sets or gets the maximum socket receive buffer in bytes. The kernel doubles this value. The minimum (doubled) value for this option is 256.
SO_RCVBUFFORCE	Using this socket option, a privileged process can perform the same task as SO_RCVBUF, but the rmem_max limit can be overridden.
SO_RCVLOWAT	Specify the minimum number of bytes in the buffer until the socket layer will pass the data to the user on receiving.
SO_SNDLOWAT	Specify the minimum number of bytes in the buffer until the socket layer will pass the data to the user on receiving.
SO_RCVTIMEO	Specify the receiving or sending timeouts until reporting an error. The argument is a struct timeval.
SO_SNDTIMEO	Specify the receiving or sending timeouts until reporting an error. The argument is a struct timeval.
SO_REUSEADDR	Indicates that the rules used in validating addresses supplied in a bind() call should allow reuse of local addresses.
SO_SNDBUF	Sets or gets the maximum socket send buffer in bytes. The kernel doubles this value when it is set.
SO_SNDBUFFORCE	Privileged process can perform the same task as SO_SNDBUF, but the wmem_max limit can be overridden.
SO_TIMESTAMP	Enable or disable the receiving of the SO_TIMESTAMP control message.

Table 49: Socket level options for `setsockopt()` as defined in `<errno.h>` [8]

B.4 IP level options for `setsockopt()`

Flag	Description
<code>IP_ADD_MEMBERSHIP</code>	Join a multicast group. Argument is an <code>ip_mreqn</code> structure.
<code>IP_ADD_SOURCE_MEMBERSHIP</code>	Join a multicast group and allow receiving data only from a specified source.
<code>IP_BLOCK_SOURCE</code>	Stop receiving multicast data from a specific source in a given group.
<code>IP_DROP_MEMBERSHIP</code>	Leave a multicast group.
<code>IP_DROP_SOURCE_MEMBERSHIP</code>	Leave a source-specific group-that.
<code>IP_FREEBIND</code>	If enabled, this boolean option allows binding to an IP address that is non-local/does not exist.
<code>IP_HDRINCL</code>	If enabled, the user supplies an IP header in front of the user data. Only valid for <code>SOCK_RAW</code> sockets.
<code>IP_MSFILTER</code>	This option provides access to the advanced full-state filtering API.
<code>IP_MTU_DISCOVER</code>	Set or receive the Path MTU Discovery setting for a socket.
<code>IP_MULTICAST_IF</code>	Set the local device for a multicast socket.
<code>IP_MULTICAST_LOOP</code>	Set or read an argument that determines if multicast packets should be looped back to the local sockets.
<code>IP_MULTICAST_TTL</code>	Set or read the time-to-live value of outgoing multicast packets for this socket.
<code>IP_NODEFRAG</code>	If enabled (nonzero), the reassembly of outgoing packets is disabled in the netfilter layer.
<code>IP_OPTIONS</code>	Set or get the IP options to be sent with every packet from this socket.
<code>IP_PKTINFO</code>	Pass an <code>IP_PKTINFO</code> ancillary message that supplies information about the incoming packet.
<code>IP_RECVERR</code>	Enable extended reliable error message passing. On a datagram socket, all generated errors are stored in a per-socket error queue.
<code>IP_RECVTOS</code>	If enabled the <code>IP_TOS</code> ancillary message is passed with incoming packets.
<code>IP_RECVTTL</code>	If set, pass a <code>IP_TTL</code> control message with the received packets TTL. Not supported for <code>SOCK_STREAM</code> sockets.
<code>IP_RETOPTS</code>	Identical to <code>IP_RECVOPTS</code> , but returns raw unprocessed options with timestamp and route record options not filled in for this hop.
<code>IP_ROUTER_ALERT</code>	Pass all to-be forwarded packets with the IP Router Alert option set to this socket. Only valid for raw sockets.
<code>IP_TOS</code>	Set or get the TOS field that is sent with every IP packet originating from this socket.
<code>IP_TRANSPARENT</code>	Setting this boolean option enables transparent proxying on this socket.
<code>IP_TTL</code>	Set or get the current time-to-live field that is used in every packet sent from this socket.
<code>IP_UNBLOCK_SOURCE</code>	Unblock previously blocked multicast source.

Table 50: IP level options for `setsockopt()` as defined in `<errno.h>` [8]

B.5 Errno flags for connect()

Flag	Description
EACCES	For UNIX domain sockets, which are identified by path-name: Write permission is denied on the socket file, or search permission is denied for one of the directories in the path prefix.
EACCES	The user tried to connect to a broadcast address without having the socket broadcast flag enabled or the connection request failed because of a local firewall rule.
EPERM	The user tried to connect to a broadcast address without having the socket broadcast flag enabled or the connection request failed because of a local firewall rule.
EADDRINUSE	Local address is already in use.
EAFNOSUPPORT	The passed address didn't have the correct address family in its source address family field.
EAGAIN	No more free local ports or insufficient entries in the routing cache.
EALREADY	The socket is non-blocking and a previous connection attempt has not yet been completed.
EBADF	The file descriptor is not a valid index in the descriptor table.
ECONNREFUSED	No-one listening on the remote address.
EFAULT	The socket structure address is outside the user's address space.
EINPROGRESS	The socket is non-blocking and the connection cannot be completed immediately.
EINTR	The system call was interrupted by a signal that was caught.
EISCONN	The socket is already connected.
ENETUNREACH	Network is unreachable.
ENOTSOCK	The file descriptor is not associated with a socket.
ETIMEDOUT	Timeout while attempting connection.

Table 51: Errno flags for connect () as defined in <errno.h>
[8]

B.6 `ioctl()` flags defined in `bpf.h`

Constant	Description
<code>BIOCGBLEN</code>	(<code>u_int</code>) Returns the required buffer length for reads on bpf files.
<code>BIOCSBLEN</code>	(<code>u_int</code>) Sets the buffer length for reads on bpf files.
<code>BIOCGDLT</code>	(<code>u_int</code>) Returns the type of the data link layer underlying the attached interface.
<code>BIOCPROMISC</code>	(<code>u_int</code>) Forces the interface into promiscuous mode.
<code>BIOCFLUSH</code>	(<code>u_int</code>) Flushes the buffer of incoming packets, and resets the statistics that are returned by <code>BIOCGSTATS</code> .
<code>BIOCGETIF</code>	(<code>struct ifreq</code>) Returns the name of the hardware interface that the file is listening on.
<code>BIOCSETIF</code>	(<code>struct ifreq</code>) Sets the hardware interface associated with the file.
<code>BIOCGRTIMEOUT</code>	(<code>struct timeval</code>) Set or get the read timeout parameter.
<code>BIOCGSTATS</code>	(<code>struct bpf_stat</code>) Returns packet statistics.
<code>BIOCIMMEDIATE</code>	(<code>u_int</code>) Enable or disable immediate mode, based on the truth value of the argument. When immediate mode is enabled, reads return immediately upon packet reception. Otherwise, a read will block until either the kernel buffer becomes full or a timeout occurs.
<code>BIOCSETFNR</code>	(<code>struct bpf_program</code>) Sets the read filter program used by the kernel to discard uninteresting packets.
<code>BIOCSETWF</code>	(<code>struct bpf_program</code>) Sets the write filter program used by the kernel to control what type of packets can be written to the interface.
<code>BIOCVERSION</code>	(<code>struct bpf_version</code>) Returns the major and minor version numbers of the filter language currently recognized by the kernel.
<code>BIOCGHRCMPLT</code>	(<code>u_int</code>) Set or get the status of the header complete flag. Set to zero if the link level source address should be filled in automatically by the interface output routine.
<code>BIOCGDIRECTION</code>	(<code>u_int</code>) Set or get the setting determining whether incoming, outgoing, or all packets on the interface should be returned by BPF (<code>BPF_D_IN</code> = only incoming, <code>BPF_D_INOUT</code> = packets originating locally and remotely, <code>BPF_D_OUT</code> = only outgoing packets).
<code>BIOCGTSTAMP</code>	(<code>u_int</code>) Set or get format and resolution of the time stamps returned by BPF. See Manpage for further details.
<code>BIOCFEEDBACK</code>	(<code>u_int</code>) Set packet feedback mode. This allows injected packets to be fed back as input to the interface when output via the interface is successful.
<code>BIOCLOCK</code>	Set the locked flag on the bpf descriptor. This prevents the execution of <code>ioctl</code> commands which could change the underlying operating parameters of the device.
<code>BIOCSETBUFMODE</code>	(<code>u_int</code>) Get or set the current bpf buffering mode.
<code>BIOCSETZBUF</code>	(<code>struct bpf_zbuf</code>) Set the current zero-copy buffer locations.
<code>BIOCGETZMAX</code>	(<code>size_t</code>) Get the largest individual zero-copy buffer size allowed.
<code>BIOCROTZBUF</code>	Force ownership of the next buffer to be assigned to userspace, if any data present in the buffer.

Table 52: `ioctl()` flags defined in `<bpf.h>` [2]