

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/339084847>

Fast Packet Processing with eBPF and XDP: Concepts, Code, Challenges, and Applications

Article in ACM Computing Surveys · February 2020

DOI: 10.1145/3371038

CITATIONS

120

READS

11,199

6 authors, including:



[Marcos Augusto M. Vieira](#)

Federal University of Minas Gerais

181 PUBLICATIONS 2,589 CITATIONS

SEE PROFILE



[Elerson Santos](#)

Federal University of Minas Gerais

17 PUBLICATIONS 174 CITATIONS

SEE PROFILE



[Eduardo Pinto](#)

Federal University of Minas Gerais

9 PUBLICATIONS 176 CITATIONS

SEE PROFILE



[Luiz Vieira](#)

Federal University of Minas Gerais

153 PUBLICATIONS 4,269 CITATIONS

SEE PROFILE

Fast Packet Processing with eBPF and XDP: Concepts, Code, Challenges, and Applications

MARCOS A. M. VIEIRA, MATHEUS S. CASTANHO, RACYUS D. G. PACÍFICO,
ELERSON R. S. SANTOS, EDUARDO P. M. CÂMARA JÚNIOR, and LUIZ F. M. VIEIRA,
Universidade Federal de Minas Gerais, Brazil

Extended Berkeley Packet Filter (eBPF) is an instruction set and an execution environment inside the Linux kernel. It enables modification, interaction, and kernel programmability at runtime. eBPF can be used to program the eXpress Data Path (XDP), a kernel network layer that processes packets closer to the NIC for fast packet processing. Developers can write programs in C or P4 languages and then compile to eBPF instructions, which can be processed by the kernel or by programmable devices (e.g., SmartNICs). Since its introduction in 2014, eBPF has been rapidly adopted by major companies such as Facebook, Cloudflare, and Netronome. Use cases include network monitoring, network traffic manipulation, load balancing, and system profiling. This work aims to present eBPF to an inexperienced audience, covering the main theoretical and fundamental aspects of eBPF and XDP, as well as introducing the reader to simple examples to give insight into the general operation and use of both technologies.

CCS Concepts: • **Networks** → **Programming interfaces; Middle boxes / network appliances**; End nodes;

Additional Key Words and Phrases: Computer networking, packet processing, network functions

ACM Reference format:

Marcos A. M. Vieira, Matheus S. Castanho, Racyus D. G. Pacífico, Elerson R. S. Santos, Eduardo P. M. Câmara Júnior, and Luiz F. M. Vieira. 2019. Fast Packet Processing with eBPF and XDP: Concepts, Code, Challenges, and Applications. *ACM Comput. Surv.* 53, 1, Article 16 (December 2019), 36 pages.
<https://doi.org/10.1145/3371038>

1 INTRODUCTION

The increase in Internet traffic and the growing complexity of services offered in data center networks have required ever-higher packet processing rates. Also, the dynamicity of service demands requires the network to adapt quickly to maintain adequate levels of quality of service and use available resources efficiently. However, computer networks have been traditionally developed in a static way, embedding the implementation of communication protocols in the hardware of network devices, making it difficult to adapt to current demands.

In recent years several proposals have been made to add more programmability to networks. Among them, we can highlight the SDN [28, 43] and NFV [48] paradigms, new computer systems

Authors' address: M. A. M. Vieira (corresponding author), M. S. Castanho, R. D. G. Pacífico, E. R. S. Santos, E. P. M. Câmara Júnior, and L. F. M. Vieira, Universidade Federal de Minas Gerais, Av. Antônio Carlos, 6627 Prédio ICEx, Belo Horizonte, MG, 31270-901, Brazil; emails: {mmvieira, matheus.castanho, racyus, elerson, epmcj, lfvieira}@dcc.ufmg.br.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2019 Association for Computing Machinery.

0360-0300/2019/12-ART16 \$15.00

<https://doi.org/10.1145/3371038>

and languages, such as POF [61], P4 [14], and more recently the extended Berkeley Packet Filter (eBPF) and the eXpress Data Path (XDP) Linux kernel network layer. This work presents eBPF and XDP assuming little to no background about the subject by the reader.

eBPF provides an instruction set and an execution environment inside the Linux kernel. It is used to modify the processing of packets in the kernel and also allows the programming of network devices. The developer writes an application in restricted C language and then compile the code into eBPF instructions. The resulting eBPF code can be processed in the kernel or by programmable devices such as SmartNICs.

XDP is the lowest layer of Linux network stack [31]. It enables developers to install programs that process packet into the Linux kernel. These programs will be called for every incoming packet. XDP is designed for fast packet processing applications while also improving programmability. In addition, it is possible to add or modify these programs without modifying the kernel source code. eBPF programs modify the (programmable) kernel operation in runtime, not requiring recompilation of the kernel.

The importance of eBPF and XDP is highlighted by its fast adoption since its introduction in the Linux kernel in 2014 by both industry and academia. Their use cases have grown rapidly to include tasks such as network monitoring, network traffic handling, load balancing, and operating system insight. Several companies already use eBPF on projects such as Facebook [26], Netronome [10], and Cloudflare [11].

We organized this tutorial as follows: the remainder of this section introduces the reader to the original BPF. In Section 1.2, the architecture of the eBPF machine is described. Section 2 presents the eBPF system. In Section 3, we describe aspects of eBPF programs, such as their structure, the types of programs available, what maps are and how to use them, the types of maps available, what helper functions are, and interaction from user space with libbpf library. In Section 4, we explain how eBPF uses hooks and present two of them: the XDP and the TC. Section 5 shows examples of eBPF programs and points the reader to extra material. In Section 6, some useful tools for developing and debugging eBPF programs are listed. Section 7 describes the existing software and hardware platforms that can process eBPF instructions. Section 8 discusses some existing industry-led research and open source projects. Section 9 presents the current limitations on eBPF and suggestions on how to overcome them. Finally, Section 10 compares eBPF with other similar technologies, and Section 11 concludes this work.

All code in this article was tested using kernel version 5.0. A stable version of the extra material is accessible on Zenodo [66]. Step-by-step instructions on how to compile, load and run each example shown throughout this text, including a VM with all tools and dependencies necessary to develop eBPF programs are available on Github [67].

1.1 BPF

Inspired by previous work on in-kernel packet filters [50], the Berkeley Packet Filter (BPF) [46] was proposed by Steven McCanne and Van Jacobson in 1992 as a solution to perform packet filtering on the kernel of Unix BSD systems. It consisted of a set of instructions and a virtual machine (VM) for executing programs written in that language.

Initially, the bytecode of an application was transferred from the user space to the kernel, where it was then checked to assure security and prevent kernel crashes. After passing the verification, the system attached the program to a socket and ran on each arriving packet. The ability to securely run programs provided by the user in the kernel proved to be a good design choice of BPF. Another highlighting factor of BPF was its simple and well-defined set of instructions. Furthermore, there existed a Just-In-Time (JIT) compilation engine for BPF in the kernel. Together, all these factors were fundamental for the good performance of the tool.

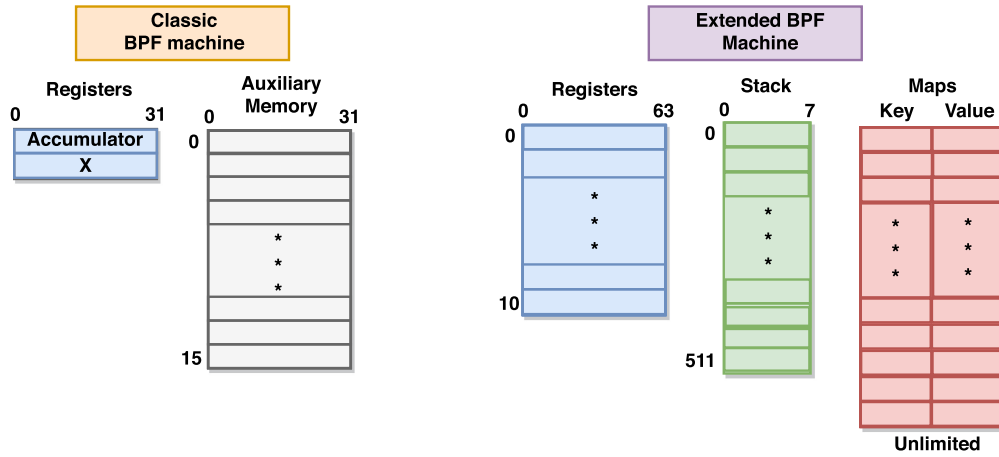


Fig. 1. BPF and eBPF processors.

Besides the bytecode instructions, BPF also defines a packet-based memory model (load instructions are implicitly made in the processed packet), two registers: accumulator (A) and index register (X), an implied program counter, and a temporary auxiliary memory. The left side of Figure 1 (Classic BPF machine) illustrates the BPF machine architecture.

The Linux kernel has supported BPF since version 2.5. There were no major changes to the BPF code until 2011, when the BPF interpreter was modified to be a dynamic translator [25]. Instead of interpreting the BPF byte code, the kernel was now able to translate BPF programs directly into x86 instructions.

One of the most prominent tools that use BPF is the libpcap library, used by the tcpdump tool. When using tcpdump to capture packets, a user can set a packet filtering expression so that only packets matching that expression are actually captured. For example, the expression “ip and tcp” captures all IPv4 packets that contain the TCP transport layer protocol. This expression can be reduced by a compiler to BPF bytecode. Code 1, based on a bpf man page [35], is a BPF program that filters packets to capture only TCP segments. The mnemonic were expanded for clarity.

Code 1 BPF program example based on [35] to only allow IPv4 TCP segments.

```

1 load 2 bytes @ [12]
2 jump equal #0x800 jump true 3 jump false 6
3 load 1 byte @ [23]
4 jump equal #6 jump true 5 jump false 6
5 return #-1
6 return #0

```

Basically, what Code 1 does is as follows:

- Instruction (1): loads two bytes (16 bits) at offset 12 of the frame into the accumulator. The offset 12 represents the packet type in the Ethernet frame.
- Instruction (2): compares the accumulator value with 0x800, which is the EtherType value for IPv4. If the result is true, then the program counter jumps (jumprtrue) to instruction (3) and otherwise jumps (jumpfalse) to instruction (6).
- Instruction (3): loads the offset 23 of the frame, as a byte, into the accumulator. The offset 23 represents the protocol field of the IPv4 packet. The count is from the beginning of the Ethernet frame.

Table 1. Description of the eBPF Register Set

Register	Description
r0	return value from functions and programs
r1 - r5	arguments passed to functions
r6 - r9	registers that are preserved during function calls
r10	stores frame pointer to access the stack

102 • Instruction (4): compares the value with the constant 6 (value of the IPv4 packet protocol
 103 field for a TCP segment). If true, then skip to instruction (5); otherwise, go to instruction
 104 (6).

105 The packet filtering program executes until it returns a result, which is usually a Boolean. Re-
 106 turning a value other than zero (instruction (5)) means that the packet has matched the filter,
 107 whereas returning zero (instruction (6)) indicates the packet does not match the filter and there-
 108 fore will be discarded.

109 1.2 Extended BPF

110 Although BPF was very useful for packet filtering, the community came to realize that other areas
 111 could also benefit from its ability to instrument the kernel. To transform it into a *universal in-
 112 kernel virtual machine* [21], a lot of improvements were introduced to both the BPF machine and
 113 its overall architecture. This new version is called eBPF (extended BPF), or simply BPF, while the
 114 original iteration became cBPF (classic BPF). eBPF was introduced in version 3.15 of the Linux
 115 kernel. The content in this section is based on the eBPF specification [60].

116 The right side of Figure 1 illustrates the eBPF engine. The number of registers has increased
 117 from 2 to 11 (of which 10 are write-registers), the registers width has changed from 32 bits to
 118 64 bits, the instruction set is now 64 bits, and the new engine has a stack of 512 bytes. Global data
 119 stores, called maps, were also included, allowing programs to persist data between executions and
 120 share information between each other and with user space. It was also added the option to call
 121 functions that run inside the kernel, called helper functions [60].

122 In cBPF, it was necessary to define the jumps for true and false cases in a program. In eBPF, it
 123 is only necessary to define the true jumps, and the false jumps follow the execution sequence of
 124 the program (called jump-fall-through).

125 The eBPF's instruction set architecture (ISA) was updated to include function calls. Those calls
 126 follow the C calling convention. Parameters are passed to functions through registers, just as it
 127 happens in native hardware. This allows mapping an eBPF function call to one hardware instruc-
 128 tion, which results in almost no overhead. eBPF uses this feature to enable helper functions, al-
 129 lowing programs to make system calls and manipulate storage (maps). The eBPF virtual machine
 130 supports dynamic loading and program reloading. This way, programs can be changed on runtime,
 131 modified, or reloaded again if necessary.

132 Table 1 describes the functionality of each eBPF register. Register r0 stores the function return
 133 value, which indicates, at the end of the computation, what action will be taken in the forwarding
 134 of the packet. Register r10 is the only read-only register, and it stores the address to the BPF stack.

135 As eBPF follows the C calling convention, arguments are passed as register values to functions.
 136 Thus registers r1-r5 are reserved for this purpose, while registers r6-r9 have their values pre-
 137 served between function calls.

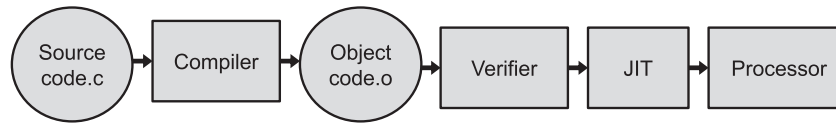


Fig. 2. eBPF Workflow.

2 EBPF SYSTEM

The eBPF system is composed of a series of components to compile, verify, and execute the source code of developed applications. This section describes more details on each of them.

2.1 Overview

The typical workflow of the eBPF system is illustrated in Figure 2. An eBPF program is written in a high-level language (mainly restricted C). The clang compiler transforms it into an ELF/object code. An ELF eBPF loader can then insert it into the kernel using a special system call. During this process, the verifier analyzes the program and upon approval the kernel performs the dynamic translation (JIT). The program can be offloaded to hardware, otherwise it is executed by the processor itself.

2.2 Compiler

Starting in version 3.7, the LLVM compiler collection has a backend for the eBPF platform. It allows the development of eBPF programs in a subset of C and generation of executable code in eBPF format through the clang compiler.

This subset of C excludes some syscalls and libraries, but it provides helper functions to manipulate eBPF maps and to perform other common tasks. Partial solutions to these restrictions are presented later in Section 9.

The main restrictions are as follows:

- eBPF can only use a subset of C language libraries. For example, the *printf()* function is not available for use;
- Non-static global variables are not allowed;
- Only bounded loops are allowed.
- Stack space is limited to 512 bytes.

There are also efforts from open source projects such as IOVisor [17, 32] and from VMWare [68] to implement a P4 [14] compiler for eBPF. Early versions already exist but are still not ready for production. Also, the BPF Compiler Collection (BCC) project [7] enables extra abstractions over the standard C code to facilitate writing and interacting with eBPF programs, as well as libbpf (Section 3.6), the main upstream library for user space interaction with eBPF.

2.3 Verifier

To ensure the integrity and security of the operating system, the kernel uses a verifier that performs static program analysis of eBPF instructions being loaded into the system. Its implementation is available at `kernel/bpf/verifier.c` in the kernel source code.

Among other things, the verifier checks whether a program is larger than allowed (current limit is 10^6 instructions), whether or not the program ends, whether the memory addresses are within the memory range allowed for the program, and how deep the execution path is. It is called after the code has been compiled and during the process of loading the program into the data plane. A good overview is presented by Miller [49].

The verifier uses two passes to decide whether to reject a program or not. In the first pass, it uses a depth-first search to check if the program instructions can be parsed into a Directed Acyclic Graph (DAG). eBPF programs that do not have backward jumps or that have only predefined size loops (which can thus have loop unroll) can be synthesized into a DAG, guaranteeing their termination. Moreover, the DAG is useful to check for unreachable instructions (the graph must have only one connected component) and to compute the worst-case execution time.

The second pass explores all possible paths from the program's first instruction. It does so by creating a state machine, where it verifies if states present correct behaviors and also keeps records of the ones it has already checked [59]. The verifier uses the states already checked for pruning and so be able to reduce its amount of work to do. It also limits the maximum length of paths to analyze. This limit was initially 64k instructions but currently equals the maximum program size allowed.

It is also worthwhile to mention two more points about the eBPF verifier. The first one is related to the fact that some eBPF functions can only be called by programs with GPL compatible licenses. Because of that, the verifier checks wither the licenses of the functions used by a program and the program's license are compatible and rejects the program if they are not.

Last, the verifier does not allow memory accesses beyond the local variables and packet boundaries to ensure the integrity and security of the kernel. To access any bytes in the packet, it is always necessary to perform a border check (as shown later in Section 5.2.1). However, each byte only needs to be checked once, unless the storage space of the packet gets modified. This way, during the analysis of the program, the verifier guarantees that all memory accesses made to the packet are in checked addresses. If the eBPF program does not do this type of check, then the verifier rejects it, and so it cannot be loaded in the kernel [31].

3 EBPF PROGRAMS

Several types of applications can be implemented with eBPF, e.g., performance analysis, packet filtering, and traffic classification, to name a few. The main advantage of the eBPF system as a whole is offering a flexible and safe programmable environment inside the Linux kernel. For example, eBPF programs can be loaded and modified during runtime and are capable of interacting with kernel elements such as kprobes, perf events, sockets, and routing tables [44].

However, the subsystems and functionalities available to an eBPF program depend on where it is loaded in the kernel, i.e., which layer or subsystem it is attached to, which is defined by a program's type. In this section, we discuss different eBPF program types, present key-value store data structures called maps, and also show some helper functions available to eBPF programs.

3.1 How and When Are eBPF Programs Executed?

To execute an eBPF program, it is first necessary to attach it to an interface that allows custom programming. This interface is called a hook. Hooks allow the registration of programs for certain events. In Section 4, we describe two Linux kernel hooks to which eBPF programs can be attached, XDP and TC.

eBPF programs execute whenever there is an event for which they were registered. In Computer Networking, common events are sending or receiving a packet.

3.2 Program Types

Each eBPF program has a type, which determines three important aspects: What is the input passed to it (its context), which helper functions it is allowed to use, and to which kernel hook it will be attached. For example, two of the many types of eBPF programs are socket filter and tracing. The input parameter for a socket filter program is a socket buffer, containing packet metadata generated

by the kernel but stripped of L2 and L3 information. A tracing program, however, receives a set of register values. Also, the subsets of helper functions available for these two types are not the same, although there is an overlap of common general-purpose functions.

Supported program types are defined on the header file `linux/bpf.h` by the `enum bpf_prog_type`. On version 5.3-rc6, the kernel offers a total of 25 valid different program types, some of which are listed below:

- **BPF_PROG_TYPE_SOCKET_FILTER**: program to perform socket filtering;
- **BPF_PROG_TYPE_SCHED_CLS**: program to perform traffic classification at the TC layer;
- **BPF_PROG_TYPE_SCHED_ACT**: program to add actions to the TC layer;
- **BPF_PROG_TYPE_XDP**: program to be attached to the eXpress Data Path hook;
- **BPF_PROG_TYPE_LWT_{IN, OUT or XMIT}**: programs for Layer-3 tunnels;
- **BPF_PROG_TYPE_SOCKET_OPS**: program to catch and set socket operations such as retransmission timeouts, passive/active connection establishment, and so on;
- **BPF_PROG_TYPE_SK_SKB**: program to access socket buffers and socket parameters (IP addresses, ports, etc) and to perform packet redirection between sockets;
- **BPF_PROG_TYPE_FLOW_DISSECTOR**: program to do flow dissection, i.e., to find important data in network packet headers.

These program types are related to networking, which is the focus of this work. However, there are other program types for kernel tracing/monitoring (e.g., `BPF_PROG_TYPE_PERF_EVENT`, `BPF_PROG_TYPE_KPROBE` and `BPF_PROG_TYPE_TRACEPOINT`), cgroups (e.g., `BPF_PROG_TYPE_CGROUP_SKB` and `BPF_PROG_TYPE_CGROUP_SOCK`) and others [44, 51]. The full list of supported program types can be obtained directly from the kernel source code with the following command:

```
$ git grep -W 'bpf_prog_type {' include/uapi/linux/bpf.h
```

3.3 Maps

Maps are generic key-value stores available to eBPF programs. Keys and values are treated as binary blobs, allowing the storage of user-defined data structures and types, whose sizes must be informed during map definition.

Maps are created using the `bpf` system call, allowing map manipulation through the map's file descriptor. This is done by passing the command `BPF_MAP_CREATE` (defined by `enum bpf_cmd`) and the `bpf_attr` union with extra parameters to the `bpf` system call:

```
bpf( BPF_MAP_CREATE, &bpf_attr, sizeof(bpf_attr) ).
```

In this case, the following attributes should be set on `bpf_attr`:

- (1) *map_type*: type of the map to be created;
- (2) *key_size*: number of bytes to store the key;
- (3) *value_size*: number of bytes to store the value;
- (4) *max_entries*: number of rows in the map.

A user-space process can create multiple maps, and they can be accessed by both user-space processes and eBPF programs loaded in the kernel, enabling data exchange between the two environments. To access a map, an eBPF program needs to declare a special global variable, of type `struct bpf_map_def` (defined by `libbpf`), in the `maps` ELF section. During the load process, the file loader uses the syscall above to create any declared maps and pass their file descriptors to the program, which are later converted into actual pointers by the verifier for use at run time. Code 2 shows an example where a `BPF_PROG_TYPE_ARRAY` map named `mapname` is declared.

Code 2 Map declaration example.

```

1  struct bpf_map_def SEC("maps") mapname = {
2      .type = BPF_MAP_TYPE_ARRAY,
3      .key_size = sizeof( uint32_t ),
4      .value_size = sizeof( long ),
5      .max_entries = 256,
6  };

```

262 **3.3.1 Map Types.** There are many different map types available for eBPF programs, and they
 263 are defined in the *enum* `bpf_map_type`, from `linux/bpf.h`. Each map type provides a different
 264 behavior, some of them being used generically, while others have specific use cases.

265 Examples of eBPF map codes are provided by the kernel. Version 5.3-rc6 of the kernel lists a
 266 total of 24 valid different map types. Some of them are

- 267 • **BPF_MAP_TYPE_ARRAY:** a map where entries are indexed by a number, as in a high-level
 268 programming language array. It follows the RAM model, where the input to query an item
 269 is an address.
- 270 • **BPF_MAP_TYPE_PROG_ARRAY:** a map that stores references to eBPF programs. Its use allows,
 271 for example, the call of subprograms to deal with specific situations.
- 272 • **BPF_MAP_TYPE_HASH:** stores entries using a hash function.
- 273 • **BPF_MAP_TYPE_PERCPU_HASH:** a map that is similar to the **BPF_MAP_TYPE_HASH**. Allows the
 274 creation of a hash table for each processor core.
- 275 • **BPF_MAP_TYPE_LRU_HASH:** a map that stores entries using hash function. When the table is
 276 full, the policy to remove element LRU, i.e., the elements to be removed are the ones that
 277 were last used the longest.
- 278 • **BPF_MAP_TYPE_LRU_PERCPU_HASH:** allows the creation of a hash table for each processor
 279 core with LRU remove policy.
- 280 • **BPF_MAP_TYPE_PERCPU_ARRAY:** a map that is similar to the **BPF_MAP_TYPE_ARRAY**. Allows
 281 the creation of an array for each processor core.
- 282 • **BPF_MAP_TYPE_LPM_TRIE:** longest-prefix match (LPM) trie.
- 283 • **BPF_MAP_TYPE_ARRAY_OF_MAPS:** an array to store references to eBPF maps.
- 284 • **BPF_MAP_TYPE_HASH_OF_MAPS:** a hash table to store references to eBPF maps.
- 285 • **BPF_MAP_TYPE_DEVMAP:** stores reading references of network devices.
- 286 • **BPF_MAP_TYPE_SOCKMAP:** stores socket references. It can be used to implement socket redi-
 287 rection, for example.
- 288 • **BPF_MAP_TYPE_QUEUE:** a map with behavior similar to that of a queue.
- 289 • **BPF_MAP_TYPE_STACK:** a map with behavior similar to that of a stack.

290 The list of all supported map types can be obtained directly from the kernel source code with
 291 the following command:

```
292 $ git grep -W 'bpf_map_type {' include/uapi/linux/bpf.h
```

293 **3.3.2 Lifetime of Maps and Map Pinning.** Every eBPF object (programs, maps, and debug info)
 294 has a reference counter (`refcnt`) that is maintained by the kernel [63]. When a user-space process
 295 creates a map with the call `bpf_create_map()`, the kernel initializes the map `refcnt` to 1. The
 296 kernel then increments the map `refcnt` whenever a new eBPF program that uses the map is loaded
 297 and decrements it whenever one of them is closed. The map `refcnt` will also be decremented
 298 when the process that created it exits (or crashes). When a `refcnt` reaches zero, a memory free is
 299 triggered, destroying the eBPF object related to the counter. This flows represents, in a simple way,
 300 the lifetime of an eBPF map.

The scheme just described allows sharing the same eBPF map between several programs at once. It keeps the map alive as long as its parent process or some eBPF program that uses it is alive. However, there is also another way to keep eBPF maps alive: by doing map pinning.

A user-space process can pin a map (or any other eBPF object) to the BPF file system, a minimal kernel space file system located at `/sys/fs/bpf/`. When a map is pinned to this file system, the kernel increments its `refcnt`, which allows it to stay alive even though no program is using it. Similarly, when a map is unpinned, its `refcnt` gets decremented, and again it may be destroyed if it is not being used.

The map pinning can be done in several ways: by using the `bpf()` system call from user space (with command `BPF_OBJ_PIN`), through `libbpf` (Section 3.6), by using `bpf tool` (Section 6.2), or by using a special map structure provided by `iproute2` (Section 6.1). This alternative structure, called `bpf_elf_map`, is compatible with the one provided by the kernel and can be used by eBPF programs instead of `bpf_map_def`. It exposes extra members, such as pinning, which can be used to define the map scope. This field can receive three distinct values: `PIN_GLOBAL_NS`, `PIN_OBJECT_NS`, and `PIN_NONE`.

Maps created with `PIN_OBJECT_NS` have local scope, being unique to the program that declared them. As a consequence, maps with the same declaration can co-exist in different programs. In this case, a specific directory will be created in the BPF file system to store the nodes corresponding to those maps. If the value `PIN_OBJECT_GLOBAL` is used, then the map is created with a global scope, enabling it to be shared by multiple programs. This map will receive an entry in the directory `globals` in the pseudo-file system. `PIN_NONE` indicates that the map should not be fixed in the file system, disabling sharing it with other applications. Finally, a map can be unpinned by removing its file from the BPF file system. This removal can be done using the syscall `unlink()`.

3.3.3 Locked Memory. eBPF maps use locked memory, which is a resource that is usually limited by many systems. Default limits may be too low, which may cause programs to be rejected at load time. To overcome this restriction, increase the locked memory limit to a sufficient one or even remove it entirely. This limit can be changed with `ulimit -l <size>`.

3.4 Helper Functions

eBPF differs from cBPF in several ways, one being the ability to allow programs to call the so-called helper functions. These are special functions offered by the kernel infrastructure to enable interaction with the context of each hook and other kernel facilities and structures, such as maps, routing tables, tunneling mechanisms, and so on.

Tasks performed by helper functions include interacting with maps, modifying packets, and printing messages to the kernel trace. Since there are many program types, and each has a specific execution context, the list of functions callable by a specific function represents a subset of all helper functions implemented by the kernel, which varies depending on the hook the program is attached to. For example, function `bpf_xdp_adjust_tail()` is used to remove the last bytes of a packet, effectively decreasing the packet's size. However, as the name indicates, it is only available at the XDP hook. The BCC project maintains a list of helper functions for each program type [8].

The helper functions available to eBPF programs are restricted to the list provided and implemented by the kernel. The addition of new helper functions can only be done through extensions to the kernel source code, since extensions through kernel modules are not allowed. New functions should follow a calling convention shared by all functions on eBPF programs, limiting the maximum number of input parameters to 5. Parameter passing is done through the use of registers `r1-r5`, requiring no interaction with the stack.

346 The number of helper functions available is large and increases constantly with new kernel
347 versions. Version 5.3-rc6 offers a total of 109 such functions. Some of these are highlighted below:

- 348 • **bpf_map_delete_elem**, **bpf_map_update_elem**, **bpf_map_lookup_elem**: used to remove,
349 install or update, and search elements from maps, respectively;
- 350 • **bpf_get_prandom_u32**: returns a 32-bit pseudo-random value;
- 351 • **bpf_l4_csum_replace**, **bpf_l3_csum_replace**: used to recalculate Layer-4 and Layer-3
352 checksums, respectively;
- 353 • **bpf_ktime_get_ns**: returns time since system boot, in nanoseconds;
- 354 • **bpf_redirect**, **bpf_redirect_map**: functions to redirect packets to other network devices.
355 The second allows specifying the device dynamically through a special redirection map;
- 356 • **bpf_skb_vlan_pop**, **bpf_skb_vlan_push**: remove/add, respectively, VLAN tags from a
357 packet;
- 358 • **bpf_getsockopt**, **bpf_setsockopt**: similar in functionality to user-space calls to
359 **getsockopt()** and **setsockopt()** to get/set socket options.
- 360 • **bpf_get_local_storage**: returns a pointer to a local storage area. Depending on the pro-
361 gram type, this area can be shared between multiple program instances running in parallel.

362 The declarations of helper functions are spread across several header files included in the di-
363 rectory `tools/testing/selftests/bpf` in the kernel source code. However, most of them are
364 in `bpf_helpers.h`. Some common operations to perform endianness conversion are declared by
365 `bpf_endian.h`, placed in the same folder. This file offers BPF-compatible versions of well-known
366 functions like `ntohs()` and `htons()`, in the form of `bpf_ntohs()` and `bpf_htons()`, for example.

367 As explained earlier, the kernel provides the implementation of these functions, and eBPF pro-
368 grams only need to be compiled against the header files containing their signature, with no need
369 for their `.c` counterpart. This can be done by passing the path to these files in the kernel source
370 code to `clang` using the `-I` flag. However, it is also possible to make a local copy of the header
371 files needed and avoid compiling against the kernel tree, making the code easier to compile and to
372 distribute.

373 **3.4.1 Tail Calls.** eBPF programs can call other program to run next, never returning to the
374 caller, via tail calls. They can be used, for example, to simplify complex programs and build dy-
375 namic chains of programs [62]. Tail calls are implemented as long jumps, and they reuse the cur-
376 rent stack frame to avoid creating a new one, leading to minimal overhead when compared to
377 function calls. The use of tail calls involves the use of (i) a specialized map, called program ar-
378 ray (**BPF_MAP_TYPE_PROG_ARRAY**), to store references of eBPF programs, and (ii) a helper function
379 (**bpf_tail_call**) to execute the tail calls. The program array can be filled by user space with key-
380 value pairs, where the values are the file descriptors of the eBPF programs. The helper function
381 receives three arguments: the context, a reference to the program array map, and a lookup key.

382 Tail calls have some limitations, however. As the chain of tail calls can form loops, the maximum
383 number of tail calls is currently limited to 32 to avoid infinite loops. Furthermore, eBPF only allows
384 programs of the same type to be tail called. The same is true for the translation type, which should
385 match the caller's (JITed or interpreted).

386 3.5 Return Codes

387 The codes returned by eBPF programs vary in meaning and value depending on the program type.
388 For example, an XDP program (Section 4.2), returns a verdict about what should be done with the
389 packet after processing (pass along, drop, redirect, etc.), which is defined by enum `xdp_action` in
390 `bpf.h`. TC return codes (Section 4.3) have a similar meaning, but use a different enumeration type.

Socket filters, however, use the return code to indicate the packet length to be passed to the stack, 391
being able to trim or even discard the packet entirely. 392

3.6 Interaction from User Space with libbpf 393

Although the kernel exposes the `bpf()` syscall to interact with the eBPF framework from user 394
space, it is a single tool serving many purposes, making it rather complex. A more user-friendly 395
API is offered by `libbpf` [39], which is a user-space library developed by the kernel community for 396
that matter. It is available under `tools/lib/bpf` on the kernel source code and is also distributed 397
in a stand-alone version on GitHub [40], which mirrors the corresponding files from the kernel. 398
To include this library, follow the steps on the README to compile the library and link it to your 399
code: 400

```
$ LIBBPF_DIR=<path-to-libbpf>/src 401
$ clang -I${LIBBPF_DIR}/root/usr/include/ -L${LIBBPF_DIR} myprog.c -lbpf 402
```

The root directory can be different based on the `DESTDIR` used during `libbpf` compilation. 403
Finally, include the library in the C code: 404

```
#include <bpf/libbpf.h> 405
```

Several examples available in the directory `tools/testing/selftests/bpf` demonstrate use 406
cases of this library and can serve as a good starting point. Also, the stand-alone version on GitHub 407
has detailed instructions on how to build a integrate `libbpf` into projects without requiring com- 408
piling against the kernel source. 409

This API includes a few direct wrappers of the `bpf()` system call and exposes several struc- 410
tures to help the interaction with the eBPF system. For example, the user can handle information 411
about maps, programs, and object files using `struct bpf_map`, `struct bpf_program`, and `struct` 412
`bpf_object`, respectively. Each of these object-like types has specific getters and setters, whose 413
names start with the name of the structure, followed by a double underscore and a declarative 414
name of the action to be performed. The following paragraphs list some of the most common 415
functions for each of these object types. 416

After compiling a `.c` file containing eBPF programs with `clang`, the object file generated will 417
contain several ELF sections corresponding to each program. A user-space program can interact 418
with such file using the `bpf_object__*` family of functions. Some examples include the following: 419

- **`bpf_object__open`**, and **`bpf_object__open_xattr`**: read an object file and returns a 420
pointer to a `struct bpf_object`. The `_xattr` version allows specifying the program type; 421
- **`bpf_object__load`**, and **`bpf_object__load_xattr`**: load the programs from a `struct` 422
`bpf_object` into the kernel. The `_xattr` version allows specifying the desired log level; 423
- **`bpf_object__pin_maps`**: allows handling pinning of all maps from an object file; 424
- **`bpf_object__for_each_program`**: macro to iterate over each program from an object file; 425
- **`bpf_object__find_program_by_title`**: returns the handle to a BPF program based on its 426
section name. 427

Some of the functions above also have their respective counterparts (*unload*, *close*, *unpin*). 428

Another useful set of functions included in the API allows handling programs separately and 429
can be used with the program iterator showed above, for example, to apply specific actions to each 430
program in a file. These functions have a `bpf_program__*` signature, and some are shown below: 431

- **`bpf_program__is_<type>`**, and **`bpf_program__set_<type>`**: getters and setters, respec- 432
tively, for a program's type. Each type has its own pair of functions, in which `<type>` is 433
replaced by the corresponding name (e.g., `sched_cls`, `sched_xdp`, etc.). 434

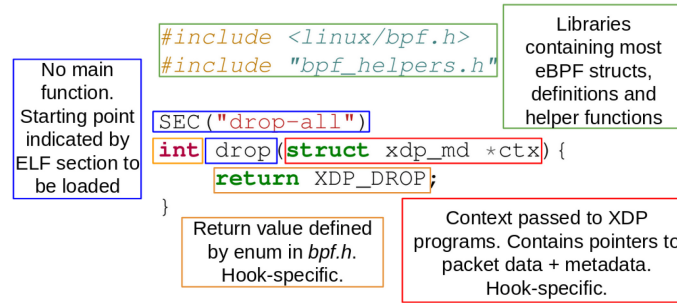


Fig. 3. Dropworld example illustrating the structure of an eBPF program.

- **bpf_program__load**: loads a given program to the kernel;
- **bpf_program__fd**: returns the file descriptor for a BPF program;
- **bpf_program__pin**: pins a given program to an specific file path;
- **bpf_program__set_ifindex**: sets a device's *ifindex* to offload maps and programs to.

At last, functions starting with `bpf_map__*` support actions on map objects, ranging from creation, information retrieval, reuse, pinning, and so on.

- **bpf_map__def**: returns basic map information (type, size, etc.);
- **bpf_map__reuse_fd**: allows the reuse of an existing map when loading a new program;
- **bpf_map__resize**: used to change the number of maximum entries allowed in a map;
- **bpf_map__fd**: returns the file descriptor for a given map;
- **bpf_map__for_each**: macro to iterate over all maps in an object file;

This section is not supposed to give an extensive discussion of all available functions, but rather give the reader a glimpse of the facilities offered by `libbpf`. The parts of the API that were left out include interaction with `perf` buffers, preprocessor helpers, and more. For the complete list of all available calls, as well as the full signature of the functions shown, please check the `libbpf.h` header file in the source code. Example code using `libbpf` will be shown later in Section 5.

3.7 Basic Program Structure

Figure 3 illustrates the basic structure of an eBPF program. It presents a simple XDP program that drops all received packets. More details about it will be given in Section 4. The library `linux/bpf.h` has all *struct* and constants definitions used by the eBPF programs, except for specific subsystems such as Traffic Control (TC) and *perf*, which need extra header files. As a rule of thumb, all eBPF programs should include this file.

Return values and the input parameter received by programs depend on the hook they will be attached. As shown in Figure 3, XDP programs receive a pointer to a `struct xdp_md`, which is explained in detail in Section 4.2.1. Programs on other hooks receive different context structures. See Section 5 for example programs on other hooks.

Note that the program shown does not contain a main function, usual on standard C programs. The program's starting point is indicated by its section in the ELF object file. When compiled, the program shown will be placed on the default `.text` section. Section 5 shows the definition of custom sections.

Below, we show the object file and disassembler output of the program example from Figure 3.

```

0:      b7 00 00 00 01 00 00 00      r0 = 1
1:      95 00 00 00 00 00 00 00      exit

```

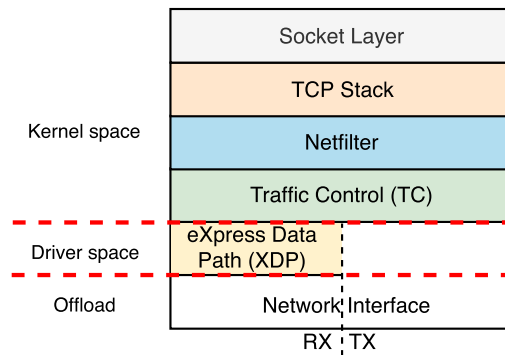



Fig. 4. Linux kernel network stack.

The first instruction writes 1 (XDP_DROP) to register `r0`. Remember from Table 1 that `r0` contains the return value of an eBPF program. The second instruction just exits. Now that the reader understands the basic program structure of an eBPF program, the next section covers XDP.

4 NETWORK HOOKS

In Computer Networking, hooks are used for intercepting packets before the call or during execution in the operating system. The Linux kernel exposes several hooks to which eBPF programs can be attached, enabling data collection and custom event handling. Although there are many hook points in the Linux kernel, we will focus on two present in the networking subsystem: eXpress Data Path (XDP) and Traffic Control (TC). Together, they can be used to process packets close to NIC on both RX and TX, enabling the development of many network applications. This section explains how eBPF can be used to program these two hooks and how programs can be loaded to each one.

4.1 Kernel's Networking Layers

Packets entering the OS are processed by several layers in the kernel, as shown in Figure 4. These layers are socket layer, TCP stack, Netfilter, Traffic Control (TC), the eXpress Data Path (XDP), and the NIC.

Packets destined to a userspace application go through all these layers and can be intercepted and modified during this process by modules such as iptables, which resides in the Netfilter layer. As explained before, eBPF programs can be attached to several places inside the kernel, enabling packet mangling and filtering.

4.2 eXpress Data Path

eXpress Data Path (XDP) is the lowest layer of the Linux kernel network stack. It is present only on the RX path, inside a device's network driver, allowing packet processing at the earliest point in the network stack, even before memory allocation is done by the OS. It exposes a hook to which eBPF programs can be attached [31].

In this hook, programs are capable of taking quick decisions about incoming packets and also performing arbitrary modifications on them, avoiding additional overhead imposed by processing inside the kernel. This renders the XDP as the best hook in terms of performance speed for applications such as mitigation of DDoS attacks.

After processing a packet, an XDP program returns an action, which represents the final verdict regarding what should be done to the packet after program exit.

Table 2. Description of XDP Action Set

Value	Action	Description
0	XDP_ABORTED	Error. Drop packet.
1	XDP_DROP	Drop packet.
2	XDP_PASS	Allow further processing by the kernel stack.
3	XDP_TX	Transmit from the interface it came from.
4	XDP_REDIRECT	Transmit packet from another interface.

498 **4.2.1 XDP Input Context.** The context seen by an XDP program is defined by the single in-
 499 put parameter passed to it by the kernel. It is of type `struct xdp_md`, defined by `bpf.h`, and is
 500 reproduced here as Code 3. Upon program execution, the `data` and `data_end` fields contain the
 501 pointers to the beginning and the end of packet data, respectively. These values must be used to
 502 guide packet access, as further explained in Section 5. The third value inside the structure is the
 503 `data_meta` pointer, which holds the address of a memory area free to be used by XDP programs
 504 to exchange packet metadata with other layers. The last two fields hold the indexes of the in-
 505 terface that received the packet and the corresponding RX queue, respectively. When accessing
 506 these two values, the BPF code is rewritten inside the kernel to access the kernel structure `struct`
 507 `xdp_rxq_info` that actually holds those values.

Code 3 Declaration of `struct xdp_md` as-is from `bpf.h`

```

struct xdp_md {
    __u32 data;
    __u32 data_end;
    __u32 data_meta;
    /* Below access go through struct xdp_rxq_info */
    __u32 ingress_ifindex; /* rxq->dev->ifindex */
    __u32 rx_queue_index; /* rxq->queue_index */
};

```

508 Although the first three fields hold pointer values, their C data type is a regular 32-byte unsigned
 509 integer. To properly use the memory addresses, a program must first cast them, which is usually
 510 done through the following code snippet, present at the beginning of almost all XDP programs:

```

    void *data_end = (void *)(long)ctx->data_end;
    void *data      = (void *)(long)ctx->data;

```

512 Here, `ctx` is the input of an XDP program, of type `struct xdp_md` shown above.

513 **4.2.2 XDP Actions.** Table 2 lists all possible XDP actions, their values, and their description.
 514 The action is specified as a program return code, which is stored at register `r0` right before the
 515 eBPF program exits.

516 The first four actions are a simple return value (no parameters), which indicate the packet should
 517 be dropped while raising an exception (`XDP_ABORTED`), dropped silently (`XDP_DROP`), passed along
 518 to the kernel stack (`XDP_PASS`) or immediately retransmitted through the same interface (`XDP_TX`).

519 The `XDP_REDIRECT` action allows an XDP program to redirect packets to (i) another NIC (physi-
 520 cal or virtual), (ii) another CPU for further processing, or (iii) an `AF_XDP` socket for userspace pro-
 521 cessing. Different from the others, this action requires a parameter to specify the redirection target.
 522 This is done through one of two helper functions: `bpf_redirect()` or `bpf_redirect_map()`. The
 523 former receives the target's interface index and is focused on network devices. The latter is a more
 524 generic alternative, which performs lookups on an auxiliary map to retrieve the final target, which

can be both net devices or CPUs. The second option is recommended, since it provides much better performance if compared to `bpf_redirect()` by batching packet transmits, and it also offers better flexibility, as map entries can be modified dynamically from user and kernel spaces.

4.2.3 XDP Modes of Operation. For increased performance, eBPF programs attached to the XDP alter the packet processing pipeline at the device driver level, which requires explicit support by the associated network driver. Some drivers for high-speed devices such as `i40e`, `nfp`, `mlx*` and the `ixgbe` family already have such functionality. On devices compatible with these drivers, XDP programs are executed directly by the driver, even before these are handled by the operating system. This is called *XDP Native* mode. BCC Project [9] maintains an up-to-date list of XDP-enabled drivers.

However, the kernel offers a compatibility mode called *XDP Generic*, which enables XDP program execution for devices without native support at the driver level. On this mode, XDP execution is done by the operating system itself, emulating native execution. This way even devices without explicit XDP support can have programs attached to them, at the cost of reduced performance due to socket buffer allocation extra steps required to perform the emulation [47].

The system automatically chooses between these two modes when loading the eBPF program. Once loaded, it is possible to check the mode of operation using the `ip` tool, as shown in the following section.

There is yet another mode of operation, *XDP Offload*. As the name suggests, the eBPF program is offloaded to compatible programmable NICs (Section 7.2.1), achieving even greater performance if compared to the other two modes. This mode should be indicated explicitly when loading the program.

4.2.4 XDP and XDP Offload Example. To demonstrate how to compile and load XDP programs, we will use the example from Figure 3. It is a simple XDP program, which drops every packet as soon as they arrive at the network interface.

After saving the example in a `dropworld.c` file, the code can be compiled into an ELF object file using the `clang` compiler:

```
$ clang -target bpf -O2 -c dropworld.c -o dropworld.o
```

The `ip` tool can load the object file into the kernel. In the example code, the program does not have any section tag, so the generated bytecode resides inside the default section (`.text`) in the ELF object file. This section should be specified when loading the program. The `-force` parameter indicates that the program should be loaded even if there is another program loaded on that interface, which will get replaced. The `[DEV]` parameter should be changed to the corresponding interface name.

```
# ip -force link set dev [DEV] xdp obj dropworld.o sec .text
```

After loading the program, the `ip` tool can also be used to verify that it is attached to the interface on the XDP hook.

```
$ ip link show dev [DEV]
```

```
DEV: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 xdp qdisc
mq state UP mode DEFAULT group default qlen 1000
    link/ether 00:16:3d:13:08:80 brd ff:ff:ff:ff:ff:ff
    prog/xdp id 27 tag f95672269956c10d jited
```

The keyword `xdp` on the first line of output indicates that an XDP program is attached to that interface in XDP Native mode. Other possible outputs could be `xdpgeneric` and `xdpoffload` for

Table 3. Description of TC Set of Actions

Value	Action	Description
0	TC_ACT_OK	Delivers the packet in the TC queue.
2	TC_ACT_SHOT	Drop packet.
-1	TC_ACT_UNSPEC	Uses standard TC action.
3	TC_ACT_PIPE	Performs the next action, if it exists.
1	TC_ACT_RECLASSIFY	Restarts the classification from the beginning.

569 the other two modes of operation. The program can also be removed from the interface by passing
 570 the parameter off:

571 # ip link set dev [DEV] xdp off

572 In this last example, the eBPF program was executed at the XDP hook by the driver, using the
 573 CPU. To offload programs, the same method as before can be used but passing the xdpoffload
 574 parameter to the ip link set command:

575 # ip -force link set dev [DEV] xdpoffload obj dropworld.o sec .text

576 As before, to remove the program just execute the following command:

577 # ip link set dev [DEV] xdpoffload off

578 4.3 Traffic Control Hook

579 Currently, although the XDP layer is well suited for many applications, it can only process ingress
 580 traffic (packets being received). To process egress traffic (transmitting packets), the closest layer
 581 to the NIC that has access to the entire Ethernet frame is the Traffic Control (TC) layer.

582 This layer is responsible for executing traffic control policies on Linux. In it, the network admin-
 583 istrator can configure different queuing disciplines (*qdisc*) for the various packet queues present
 584 in the system, as well as add filters to deny or modify packets.

585 The TC has a special queuing discipline type called *clsact*. It exposes a hook that allows queue
 586 processing actions to be defined by eBPF programs. Pointers to the packet to be processed are
 587 delivered to the configured eBPF program as part of its input context: a struct `__sk_buff`. This
 588 structure is a UAPI for certain fields that the program is allowed to access from the kernel's socket
 589 buffer internal data structure. It has the same data and data_end pointers as struct `xdp_md` but
 590 also has much more information if compared to the XDP case. This is explained by the fact that at
 591 the TC level, the kernel has already parsed the packet to extract protocol metadata, hence the richer
 592 context information passed to the eBPF program. The entire declaration of struct `__sk_buff` is
 593 omitted for brevity but can be seen on `include/uapi/linux/bpf.h`.

594 During program execution, the input packet can be modified, and the return value indicates to
 595 TC what action should be taken for it. The library `linux/pkt_cls.h` defines the available return
 596 values. The most common ones are listed in Table 3.

597 The loading of programs on the TC hook is done using the *tc* tool, available in the *iproute2*
 598 package. The following command illustrates how to create the *clsact* *qdisc* and load an eBPF
 599 program to process the packets on interface `eth0`:

600 # tc qdisc add dev eth0 clsact

601 # tc filter add dev eth0 <direction> bpf da obj <ebpf-obj> sec <section>

602 The <direction> parameter indicates which direction the program should be associated with,
 603 which can be ingress or egress. <ebpf-obj> and <section> should be the names of the file
 604 containing the compiled eBPF code and the section to load the program, respectively.

To check whether there is any program already loaded on `eth0`, use the following command: 605
`# tc filter show dev eth0 <direction>` 606

For an example of a functional eBPF program for the TC layer and its interaction with the XDP 607
 layer, please check our repository on GitHub [67]. 608

4.4 Comparison between XDP and TC 609

Both hooks can be used for similar applications, such as DDoS mitigation, tunneling, and han- 610
 dling link layer information. However, since XDP runs before any socket buffer allocation takes 611
 place, it can reach higher throughput values than programs on TC. The latter, however, can benefit 612
 from extra parsed data available through `struct __sk_buff` and execute eBPF programs for both 613
 ingress and egress traffic, being the lowest layer on TX. 614

5 EXAMPLES 615

This section presents a few examples with in-depth code explanations to help the reader become 616
 familiar with eBPF programs. The first example describes a program that allows only IPv4 TCP 617
 segments, similarly to the BPF example in Code 1. In the second example, we show the interaction 618
 between user and kernel spaces through `libbpf`, while the third one shows how programs on the 619
 XDP and TC layers can work together to collect statistics. Finally, we point the reader to a few 620
 more external examples. 621

5.1 TCP Filter 622

The first example is a program to only accept packets with TCP segments (Code 4). This is similar 623
 to the example in Code 1, but it uses eBPF to drop packets with no TCP segment. Here we present it 624
 in two perspectives: the higher level C code and the the actual eBPF assembly-like code generated 625
 after compilation. 626

5.1.1 C Code. This program was designed to be loaded into the XDP hook, so the input param- 627
 eter of the function must be of type `struct xdp_md`, as discussed in Section 4.2.1. The bytes of 628
 the packet being processed are delimited by the `data` and `data_end` pointers, which must be used 629
 throughout the program to access the packet. Type conversions for these two values are standard, 630
 so Lines 9 and 10 should be used at the beginning of every eBPF program that accesses packet 631
 data. By using `data`, the parsing of headers can be done with the standard header files provided 632
 by Linux. 633

The main difference, however, to other common packet parsing Linux programs is that bound 634
 checks are necessary before actually accessing protocol header data. Since the kernel verifier 635
 performs strict memory bound checks (Section 2.3), every access to packet data needs to be 636
 covered by an if statement with a border check (Lines 14 and 21). Each byte only needs to be 637
 checked once, unless helper functions that modify the storage space of the packet are used (e.g., 638
`bpf_xdp_adjust_head()`). In that case, it is necessary to redo all checking after calling such func- 639
 tions. If an eBPF program does not perform this type of check, then it gets rejected during load 640
 time by the verifier and is not loaded into the kernel. 641

After extracting the packet's protocol number, the program checks if it corresponds to the TCP 642
 protocol and allows it to go through the stack (Line 26). If not, then the packet is just dropped 643
 (Line 28). 644

5.1.2 eBPF Bytecode. Upon compilation, `clang` generates an object file with the eBPF instruc- 645
 tions, which can then be loaded into the kernel. As explained in Section 2.3, the verifier creates a 646
 DAG based on the program. Figure 5 shows the respective DAG in this example. Each DAG node 647

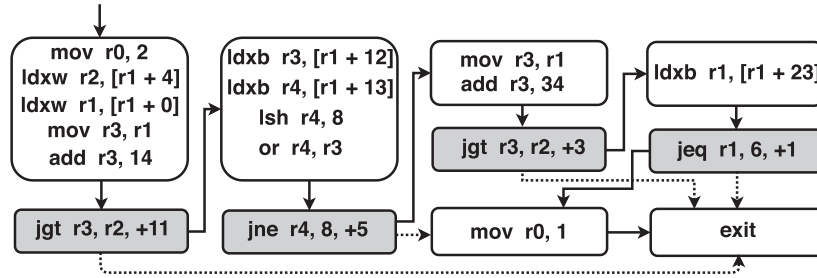


Fig. 5. Directed acyclic graph example of Code 4.

contains one or more eBPF instructions. The conditional jump nodes are the nodes that contain two output lines and a light gray background. The solid line indicates the next Acyclic Control Flow Graph (ACFG) node. Dotted lines indicate jumps to another ACFG node. In our example, there are three different kinds of conditional jump instructions: *jgt* (jump if greater), *jne* (jump if not equal), and *jeq* (jump if equal). The last number on each of this instructions indicates how many instructions to jump when the condition is valid.

Code 4 Example of a C code that checks if the packet contains an IPv4 TCP segment.

```

1  #include <linux/bpf.h>
2  #include <linux/if_ether.h>
3  #include <linux/ip.h>
4  #include <linux/tcp.h>
5  #include <linux/in.h>
6  #include "bpf_endian.h"
7
8  int isTCP( struct xdp_md *ctx ) {
9      void *data_end = (void *) (long) ctx->data_end;
10     void *data_begin = (void *) (long) ctx->data;
11     struct ethhdr* eth = data_begin;
12
13     // Check packet's size
14     if(eth + 1 > data_end)
15         return XDP_PASS;
16
17     // Check if Ethernet frame has IPv4 packet
18     if (eth->h_proto == bpf_htons( ETH_P_IP )) {
19         struct iphdr *ipv4 = (struct iphdr *) ( ((void*)eth) + ETH_HLEN );
20
21         if(ipv4 + 1 > data_end)
22             return XDP_PASS;
23
24         // Check if IPv4 packet contains a TCP segment
25         if (ipv4->protocol == IPPROTO_TCP)
26             return XDP_PASS;
27     }
28     return XDP_DROP;
29 }

```

To better understand this eBPF program, remember that register *r1* starts with a pointer to the input context stored in main memory and register *r0* stores the return value (Table 1).

In the first node, the first assembler instruction sets *r0* (from Table 1) to 2 (XDP_PASS as seen in Table 2. Moreover, load instructions compute the references to the packet's beginning and ending

from the input passed (Lines 10–12). Then the bounds of the Ethernet header are checked to guarantee valid memory access later (required by the verifier) (Lines 14–16). If the check fails, then the first jump instruction switches the flow to the *exit* instruction, ending the program. Otherwise, bytes 12 and 13 of the Ethernet header are loaded (counting from 0), which is the Ethernet type field. Then, a byte swap is done to set the endianness (Line 19).

The second jump instruction compares whether the Ethernet type is `0x0800` (Line 19). Then the bounds of the IP header are checked to guarantee valid memory access later (again, required by the verifier) (Lines 21–23). Next, the IP protocol field is extracted from the IPv4 packet, and then the program checks whether it is the TCP protocol (value 6) (Line 26) (last gray background node in the DAG). Finally, the packet is passed along to the kernel, since register (`r0`) was previously loaded with the value 2 (`XDP_PASS`), which indicates acceptance. If the IP packet does not contain a TCP header, then the value 1 (`XDP_DROP`) is loaded to `r0`. The last instruction tells that the code terminates.

5.2 User and Kernel Space Interaction

Following, we present the `xdp1` example, extracted directly from the kernel source code, present in the `samples/bpf` directory. It is divided into two parts: `xdp1_kern.c` is the actual eBPF program to be compiled and loaded into the kernel, while `xdp1_user.c` is a user-space counterpart to load the eBPF program into the kernel and interact with it through maps. We show the code for each one and discuss them separately below.

5.2.1 Kernel Space. The file `xdp1_kern.c` [54] contains an eBPF program that processes each packet on the XDP hook, extracts the corresponding IP protocol number, counts the number of packets received per protocol using a per-CPU array map called `rxcnt` and finally drops all incoming traffic.

```

1  /* Copyright (c) 2016 PLUMgrid
2  *
3  * This program is free software; you can redistribute it and/or
4  * modify it under the terms of version 2 of the GNU General Public
5  * License as published by the Free Software Foundation.
6  */
7
8  #include <uapi/linux/bpf.h>
9  #include <linux/in.h>
10 #include <linux/if_ether.h>
11 #include <linux/if_packet.h>
12 #include <linux/if_vlan.h>
13 #include <linux/ip.h>
14 #include <linux/ipv6.h>
15 #include "bpf_helpers.h"
16
17 struct bpf_map_def SEC("maps") rxcnt = {
18     .type = BPF_MAP_TYPE_PERCPU_ARRAY,
19     .key_size = sizeof(u32),
20     .value_size = sizeof(long),
21     .max_entries = 256,
22 };
23
24 static int parse_ipv4(void *data, u64 nh_off, void *data_end) {
25     struct iphdr *iph = data + nh_off;
26
27     if (iph + 1 > data_end)

```


16:20

M. A. M. Vieira et al.

```

28     return 0;
29     return iph->protocol;
30 }
31
32 static int parse_ipv6(void *data, u64 nh_off, void *data_end) {
33     struct ipv6hdr *ip6h = data + nh_off;
34
35     if (ip6h + 1 > data_end)
36         return 0;
37     return ip6h->nexthdr;
38 }
39
40 SEC("xdp1")
41 int xdp_prog1(struct xdp_md *ctx) {
42     void *data_end = (void *) (long) ctx->data_end;
43     void *data = (void *) (long) ctx->data;
44     struct ethhdr *eth = data;
45     int rc = XDP_DROP;
46     long *value;
47     u16 h_proto;
48     u64 nh_off;
49     u32 ipproto;
50
51     nh_off = sizeof(*eth);
52     if (data + nh_off > data_end)
53         return rc;
54
55     h_proto = eth->h_proto;
56
57     if (h_proto == htons(ETH_P_8021Q) || h_proto == htons(ETH_P_8021AD)) {
58         struct vlan_hdr *vhdr;
59
60         vhdr = data + nh_off;
61         nh_off += sizeof(struct vlan_hdr);
62         if (data + nh_off > data_end)
63             return rc;
64         h_proto = vhdr->h_vlan_encapsulated_proto;
65     }
66     if (h_proto == htons(ETH_P_8021Q) || h_proto == htons(ETH_P_8021AD)) {
67         struct vlan_hdr *vhdr;
68
69         vhdr = data + nh_off;
70         nh_off += sizeof(struct vlan_hdr);
71         if (data + nh_off > data_end)
72             return rc;
73         h_proto = vhdr->h_vlan_encapsulated_proto;
74     }
75
76     if (h_proto == htons(ETH_P_IP))
77         ipproto = parse_ipv4(data, nh_off, data_end);
78     else if (h_proto == htons(ETH_P_IPV6))
79         ipproto = parse_ipv6(data, nh_off, data_end);
80     else
81         ipproto = 0;
82
83     value = bpf_map_lookup_elem(&rxcnt, &ipproto);
84     if (value)

```

Fast Packet Processing with eBPF and XDP

16:21

```

85         *value += 1;
86
87         return rc;
88     }
89     char _license[] SEC("license") = "GPL";

```

The first thing to point out is that a C source file can contain many eBPF programs. They are separated into unique sections in the ELF file generated by the compiler. The section label above the corresponding function (Line 40) indicates to the compiler the name of the ELF section that will contain the program in the generated object file. This information is required while loading the code into the kernel so that the system knows which ELF section to load. Section labels are also used during map (Line 17) and program license (Line 89) declarations, both of which have fixed values. The verifier uses the license section to determine which helper functions will be available to the user, as some of them are restricted to programs declaring GPL compatible licenses.

Similarly to the previous example, the program parses the packet headers up to the IPv4 or IPv6 headers. After determining the Layer-4 protocol type (Lines 77 and 79), the program retrieves the counter for the corresponding protocol using the lookup helper function (Line 83). This function returns a pointer to the current value stored in the map if it exists, or NULL otherwise. This address can be used to change the stored data directly, without the need for a map update operation. Finally, the program returns the action that must be taken by the XDP hook for the current packet, which in this case is always XDP_DROP, indicating that the packet should be discarded.

5.2.2 User Space. The statistics collected by the kernel and stored on the map `rxcnt` are then queried by a user-space application, implemented by the file `xdp1_user.c` [55]. For brevity, we highlight below only the meaningful parts of this program.

First, instead of using the facilities of `iproute2` as shown in Section 4.2.4 to load the program to the kernel, this example implements a custom loader based on `libbpf` (Section 3.6). This method yields a higher degree of control over how the program in `xdp1_kern.c` is loaded to the kernel, which can be modified programmatically from the user-space application. For such, `libbpf.h` and `bpf.h` are included to allow interaction with the eBPF system from user space:

```

1 // SPDX-License-Identifier: GPL-2.0-only
2 /* Copyright (c) 2016 PLUMgrid
3 */
21 #include "bpf/bpf.h"
22 #include "bpf/libbpf.h"

```

The program information to be loaded is passed through the `bpf_prog_load_attr` structure, including the program type, the object file containing the program, and the interface identifier to which it should be associated.

```

73 struct bpf_prog_load_attr prog_load_attr = {
74     .prog_type = BPF_PROG_TYPE_XDP,
75 };
112 snprintf(filename, sizeof(filename), "%s_kern.o", argv[0]);
113 prog_load_attr.file = filename;

```

This structure is then used to load the program into the XDP hook. In the case of success, after the call, the variables `obj` and `prog_fd` contain the detailed information of the code already loaded and its file descriptor, respectively. The descriptor is used to identify the program from the others currently loaded in the kernel, which is necessary for future interactions with this program.

```

115 if (bpf_prog_load_xattr(&prog_load_attr, &obj, &prog_fd))
116     return 1;

```

16:22

M. A. M. Vieira et al.

711 After loading the eBPF program into the kernel, the reference to the `rxcnt` map is obtained. The
 712 function `bpf_map__next` returns an iterator for the list of maps declared in the program. Since in
 713 this case there is only one declared map, the value of that iterator can be used to obtain the file
 714 descriptor referring to it. The `libbpf.h` library also offers other functions to get map descriptors
 715 by name or by index in the map list.

```

118 map = bpf_map__next(NULL, obj);
119 if (!map) {
120     printf("finding a map in obj file failed\n");
121     return 1;
122 }
123 map_fd = bpf_map__fd(map);

```

716 Finally, the eBPF program is ready to be attached to an interface, and the userspace application
 717 can enter the infinite loop inside the `poll_stats()` function.

```

130 if (bpf_set_link_xdp_fd(ifindex, prog_fd, xdp_flags) < 0) {
131     printf("link set xdp fd failed\n");
132     return 1;
133 }
134
135 poll_stats(map_fd, 2);

```

718 The poll function `poll_stats()`, using the file descriptor of map `rxcnt`, performs periodic
 719 lookups to it and lists all the existing entries along with the statistics calculated so far.

```

35 static void poll_stats(int map_fd, int interval)
36 {
37     unsigned int nr_cpus = bpf_num_possible_cpus();
38     __u64 values[nr_cpus], prev[UINT8_MAX] = { 0 };
39     int i;
40
41     while (1) {
42         __u32 key = UINT32_MAX;
43
44         sleep(interval);
45
46         while (bpf_map_get_next_key(map_fd, &key, &key) != -1) {
47             __u64 sum = 0;
48
49             assert(bpf_map_lookup_elem(map_fd, &key, values) == 0);
50             for (i = 0; i < nr_cpus; i++)
51                 sum += values[i];
52             if (sum > prev[key])
53                 printf("proto %u: %10llu pkt/s\n",
54                     key, (sum - prev[key]) / interval);
55             prev[key] = sum;
56         }
57     }
58 }

```

720 Note that since `rxcnt` is a per-cpu array map, stored data are actually spread across multiple
 721 CPUs. The helper `bpf_num_possible_cpus` (Line 37) retrieves the number of CPUs used by the
 722 program, which is used to set the size of the array that will hold to data from each CPU (`values`).
 723 Then, an infinite loop queries the entire map from time to time to retrieve recently collected statis-
 724 tics. This is done with the `bpf_map_get_next_key` iterator function, which yields one map key at
 725 a time, enabling iterating through all entries in the map in order. Remember from the kernel-side

program that the keys are IP protocols numbers, so given the map size declared (256), iterating through all keys corresponds to iterating through all IP protocol numbers possible.

The user-space version of `bpf_map_lookup_elem` is used (Line 49) to actually read the map values associated with the corresponding key from all CPUs at the same time, which are then stored by the helper function in the values array. These values are then added to get the overall statistic (Line 51), and if the value obtained is greater than what was seen the last time, then the difference is printed to standard output, showing the user how many packets with that specific protocol number were received during the sampling interval.

This example shows how user and kernel spaces can interact through eBPF programs and maps. All data collection and packet handling on the fast path (kernel) is executed by a minimal, optimized eBPF program, while an agent retrieves the data periodically on the slow path (user space) and can take actions based on it, for example, display to the user. This is a very powerful and useful approach that can be applied and extended to many different scenarios.

5.3 Cooperation between XDP and TC

The following example consists of two separate eBPF programs, one to be attached to the XDP layer and another to TC. Together, they track the number of packets and bytes exchanged between two different IPv4 addresses and store these pieces of information on a map, tracking both RX and TX. This example demonstrates some of the unique facilities offered by `iproute2` (Section 6.1), how programs can be loaded without the need of a custom user-space program, and how programs in different layers can interact through maps.

```

1  #include <stdbool.h>
2  #include <stdint.h>
3  #include <stdlib.h>
4  #include <linux/in.h>
5  #include <linux/bpf.h>
6  #include <linux/ip.h>
7  #include <linux/tcp.h>
8  #include <linux/if_ether.h>
9  #include <linux/pkt_cls.h>
10 #include <iproute2/bpf_elf.h>
11
12 #include "bpf_endian.h"
13 #include "bpf_helpers.h"
14
15 struct pair {
16     uint32_t lip; // local IP
17     uint32_t rip; // remote IP
18 };
19
20 struct stats {
21     uint64_t tx_cnt;
22     uint64_t rx_cnt;
23     uint64_t tx_bytes;
24     uint64_t rx_bytes;
25 };
26
27 struct bpf_elf_map SEC("maps") trackers = {
28     .type = BPF_MAP_TYPE_HASH,
29     .size_key = sizeof(struct pair),
30     .size_value = sizeof(struct stats),
31     .max_elem = 2048,
32     .pinning = 2, // PIN_GLOBAL_NS

```

16:24

M. A. M. Vieira et al.

```

33  };
34
35  static bool parse_ipv4(bool is_rx, void* data, void* data_end, struct pair *pair){
36      struct ethhdr *eth = data;
37      struct iphdr *ip;
38
39      if(data + sizeof(struct ethhdr) > data_end)
40          return false;
41
42      if(bpf_ntohs(eth->h_proto) != ETH_P_IP)
43          return false;
44
45      ip = data + sizeof(struct ethhdr);
46
47      if ((void*) ip + sizeof(struct iphdr) > data_end)
48          return false;
49
50      pair->lip = is_rx ? ip->daddr : ip->saddr;
51      pair->rip = is_rx ? ip->saddr : ip->daddr;
52
53      return true;
54  }
55
56  static void update_stats(bool is_rx, struct pair *key, long long bytes){
57      struct stats *stats, newstats = {0,0,0,0};
58
59      stats = bpf_map_lookup_elem(&trackers, key);
60      if(stats){
61          if(is_rx){
62              stats->rx_cnt++;
63              stats->rx_bytes += bytes;
64          }else{
65              stats->tx_cnt++;
66              stats->tx_bytes += bytes;
67          }
68      }else{
69          if(is_rx){
70              newstats.rx_cnt = 1;
71              newstats.rx_bytes = bytes;
72          }else{
73              newstats.tx_cnt = 1;
74              newstats.tx_bytes = bytes;
75          }
76
77          bpf_map_update_elem(&trackers, key, &newstats, BPF_NOEXIST);
78      }
79  }
80
81  SEC("rx")
82  int track_rx(struct xdp_md *ctx)
83  {
84      void *data_end = (void *) (long) ctx->data_end;
85      void *data = (void *) (long) ctx->data;
86      struct pair pair;

```

Fast Packet Processing with eBPF and XDP

16:25

```

87
88     if(!parse_ipv4(true,data,data_end,&pair))
89         return XDP_PASS;
90
91     // Update RX statistics
92     update_stats(true,&pair,data_end-data);
93
94     return XDP_PASS;
95 }
96
97 SEC("tx")
98 int track_tx(struct __sk_buff *skb)
99 {
100     void *data_end = (void *) (long)skb->data_end;
101     void *data = (void *) (long)skb->data;
102     struct pair pair;
103
104     if(!parse_ipv4(false,data,data_end,&pair))
105         return TC_ACT_OK;
106
107     // Update TX statistics
108     update_stats(false,&pair,data_end-data);
109
110     return TC_ACT_OK;
111 }

```

Besides standard C types, maps can also handle user-defined structures. For example, struct pair and struct stats (Lines 15 and 20) are used here as key and value, respectively, for the trackers map (Line 27), where communication statistics will be stored. The statistics consist on the number of RX and TX packets and bytes, which are tracked for each unique pair of IPv4 addresses seen by an interface.

Note that the map definition is different from the one in the previous example. The bpf_elf_map structure is used by iproute2, hence the inclusion of the header file iproute2/bpf_elf.h, which is added to the system when iproute2 is installed from source. Besides slightly different field names, it also contains an extra pinning field, which can be used to define the map's scope, as discussed in Section 3.3.2. The value set for this field determines there should be a single trackers map, which will be shared by both programs, instead of a separate copy for each.

This same source file holds the two programs use, in two different ELF sections. Section rx (Line 81) has an XDP program that will process packets as soon as they arrive on the system. The function parse_ipv4 (Line 35) executes all parsing necessary to extract source and IP addresses from the packet. Besides the pointers to packet data and a buffer to place the address information in, it also receives a Boolean value as an input parameter. It indicates if this is an incoming or outgoing packet. Since we want to track each pair of communicating IPs, packets flowing in both directions should be tracked by the same entry in the map. Thus, the program interprets source and destination IP addresses as local or remote. On RX, the destination IP is local (the host's or possibly a guest), whereas the source IP is remote. On TX, it is the other way around.

Once the IP pair is extracted from the packet, it is passed to the function update_stats to update the trackers map accordingly. It also receives a parameter is_rx, which will indicate if the packet seen is incoming or outgoing. The function retrieves the previous statistics for the corresponding address pair (Line 59) and checks if an entry already exists. If not, then a new entry is created with a call to bpf_map_update_elem with the BPF_NOEXIST flag (Line 77). In this case, the return value of bpf_map_update_elem is ignored, since there is not much the program can do if the call

772 fails. If stats is not null, then the values are updated directly using the that pointer (Lines 61–67).
 773 The program terminates returning a XDP_PASS code, to just pass the packet along to the stack
 774 normally.

775 Since there is no XDP layer on TX, the closest we can get to handling the packet just before
 776 sending it to the NIC is by attaching an eBPF program to the TC layer. The program on section
 777 tx does just that. It is in essence the same as the XDP version, loaded to RX, with just small
 778 changes to handle packets going on the other direction. Note that besides the arguments passed
 779 to the auxiliary functions, the input context and the return codes are different (as explained in
 780 Section 4.3), since this program will be loaded on TC.

781 In this example, two programs on different layers work independently to fill the same map. This
 782 same idea can be extended to enable sharing information between programs inside the kernel, in
 783 addition to user space.

784 5.4 Other Examples

785 The examples shown above demonstrate some basic aspects that should be considered during the
 786 development of eBPF programs. They also show how to interact with user space programs.

787 **Linux kernel:** other useful examples are offered by the source code of the Linux kernel, located
 788 in two separate directories: samples/bpf and tools/testing/selftests/bpf. Both contain pro-
 789 grams that demonstrate the use of various features and hooks on kernel stack, with new programs
 790 being added to each new version of the kernel. Most of the examples in the first directory are di-
 791 vided into two separate files, one with user space code to load the program (files ending in user . c)
 792 and the other with the implementation of the eBPF program to be loaded in the kernel (files ending
 793 in kern . c). In general, examples in this folder represent stand-alone projects that may be useful for
 794 various tasks. The examples on the second directory are used as the basis for running functional
 795 tests during the development of the kernel. The actual eBPF programs that go into the kernel
 796 are placed inside the progs/ sub-directory, while the remaining files on the root correspond to
 797 user-space programs and scripts used to load them.

798 **XDP Project:** beyond the examples in the kernel, the official XDP tutorial [70] also presents
 799 examples with step-by-step instructions, explaining different features of the XDP hook in detail.

800 **L4 load-balancer:** Netronome provides code for an XDP program called 141b that implements
 801 an L4 load balancer [53]. The program processes incoming network packets and calculates a hash
 802 value based on the source IP address along with TCP or UDP ports, to ensure that all packets from
 803 the same flow are processed in the same server. The generated hash is used as a key in an eBPF
 804 map. This eBPF map is populated with the address of the available servers for which the program
 805 can redirect the packets. This program extends and inserts an external IP header with data from
 806 the map. Then, the packet is forwarded to the corresponding server. This program could also be
 807 offloaded to a SmartNIC, saving CPU cycles.

808 **Programmable Receive Side Scaling:** Receive Side Scaling (RSS) allows to map a receiving
 809 packet to be processed by a CPU core. This is important to distribute network traffic across mul-
 810 tiple CPUs in multiprocessor systems. RSS techniques are used by many network adapters to dis-
 811 tribute the computation of packets to a set of distinct CPUs through the use of multiple queues.
 812 However, the implementations are usually proprietary and hardware-based, allowing little or no
 813 programmability. Using an eBPF program, packet distribution can be modified on demand through
 814 map values or full replacement of the loaded eBPF program [53].

815 6 TOOLS

816 This section introduces some tools that can be useful for developing and debugging eBPF programs.

6.1 iproute2

iproute2 is a set of user space tools to control, configure, and monitor the kernel's network. `ip` and `tc` are example of these tools. Both offer alternative ways of loading eBPF codes into the kernel without the need for a user space program making use of the `libbpf` library or `bpf` system call. Examples on how to use both `ip` and `tc` tools were shown previously in this text.

In addition, `iproute2` has its own interface to interact with the eBPF system, offering additional functionalities such as the ability to specify the scope of eBPF maps allocation. Unfortunately, its extra features makes it incompatible with `libbpf` loader. For example, when using `iproute2`, maps are declared using an alternative structure (`bpf_elf_map`), defined by `iproute2/bpf_elf.h`. An example usage is as follows:

```
#include <linux/bpf.h>
#include <iproute2/bpf_elf.h>

struct bpf_elf_map SEC("maps") src_mac = {
    .type = BPF_MAP_TYPE_HASH,
    .size_key = 1,
    .size_value = 6,
    .max_elem = 1,
    .pinning = PIN_GLOBAL_NS,
};
```

This structure is similar to `bpf_map_def` of `libbpf` but has extra fields, such as the pinning, used to define the map's scope (which was described in Section 3.3.2). So, `iproute2` can load `libbpf` programs, but the opposite is not true as `libbpf` does not know how to handle the extra functionalities [2].

6.2 bpftool

The `bpftool` [15] is a user-space debug utility that can also load eBPF programs, create and manipulate maps, and collect information about eBPF programs and maps. It is part of the Linux kernel tree and is available at `tools/bpf/bpftool`. Here we present just some of its functionalities.

- It can list loaded programs with the command:
`bpftool prog show`
- To print the instructions of an specific program, use:
`bpftool prog dump xlated id <id>`
- It is possible to list and print the contents of maps at runtime:
`bpftool map`
`bpftool map dump id <map id>`
- It can also perform some management operations, including loading programs, performing searches, or updating map values. An example for the last item is
`bpftool map update id 1234 key 0x01 0x00 0x00 0x00 value 0x12 0x34`
↪ `0x56 0x67`

6.3 llvm-objdump

The `llvm-objdump` tool (version 4.0 or higher) provides a disassembler to transform the compiled byte code into a format readable for humans before the user attempts to inject it into the kernel. Also, it is useful for inspecting the ELF sections of the compiled eBPF file. The command below shows how to use the disassembler:

```
$ llvm-objdump -S dropworld.o
```

852 6.4 BPF Compiler Collection

853 The open-source project BPF Compiler Collection (BCC) [7] aims to facilitate the development of
 854 eBPF programs. It provides a set of frontends that can be used to interact with the eBPF system
 855 using high-level languages such as Python and Go.

856 The project also has a series of example tools built using BCC capable of performing various
 857 tasks in the operating system. These tools can perform tasks such as analyzing the number of
 858 system calls by an application, and the time elapsed during a disc read. As they are based on eBPF
 859 programs, they can be used to analyze actual production systems with low additional overhead.

860 7 PLATFORMS

861 There are several platforms that make use of the eBPF instruction set to add programmability to
 862 different environments beyond the Linux kernel. In this section, we present the most prominent,
 863 divided by their nature: software or hardware.

864 7.1 Software

865 *7.1.1 Linux Kernel.* As discussed previously, the Linux kernel was the origin of eBPF and is
 866 the most popular platform, with the most active development being the main focus of this text.
 867 The eBPF applications inside the kernel are not only restricted to network processing but can
 868 also be applied to kernel instrumentation and monitoring. Thus, eBPF programs today represent
 869 an import toolset to perform Linux introspection and performance analysis, with important open
 870 source projects such as IOVisor [32] offering many solutions on this front.

871 *7.1.2 Userspace BPF.* The *Userspace BPF (uBPF)* [65] is an open-source project that adapts the
 872 eBPF processor to run in userspace. *uBPF* project has a copy of the eBPF interpreter and JIT com-
 873 piler stripped of all kernel-specific data structures, allowing users to embed an eBPF machine into
 874 other projects running in userspace. By leveraging the *uBPF* source code, one can easily add pro-
 875 grammability with eBPF to other tools and environments. Similarly to *uBPF*, project *rbpf* offers an
 876 alternative implementation of the eBPF VM in Rust [52]. Since *uBPF* runs on userspace, it does not
 877 have the restrictions imposed by the eBPF verifier like requiring unrolling loop.

878 However, differently from the eBPF system in the kernel, *uBPF* does not offer native support to
 879 maps and does not have any helper functions implemented. Nonetheless, it can be easily extended
 880 to support these features, as was done in Reference [34], who used it as part of the BPFabric
 881 programmable virtual switch, discussed next.

882 *7.1.3 BPFabric.* To deal with the limitations of OpenFlow, [34] propose a new SDN architecture
 883 called BPFabric. BPFabric is a switch architecture that allows the data plane to be programmed with
 884 eBPF instructions by using a modified version of *uBPF*.

885 Since eBPF programs can be dynamically modified, it allows the parsing of arbitrary protocols
 886 and the use of eBPF maps to store states. New helper functions can be developed to add support to
 887 new features and to provide services such as telemetry, statistics collection, and packet tracking.

888 The control plane hosts an agent that communicates with the data plane through the South-
 889 bound API. The agent is responsible for (i) changing the behavior of the switch, (ii) receiving
 890 packets, (iii) reporting events, and (iv) reading and updating table entries. When the agent re-
 891 ceives the compiled code from the controller, it uses the eBPF ELF Loader to modify the switch
 892 pipeline. The eBPF Loader calls the verifier to check the code. On success, it must perform the
 893 allocation of the required eBPF tables and convert the received byte code into a switch-specific
 894 format.

Fast Packet Processing with eBPF and XDP

16:29

When receiving a packet, the eBPF processor executes the previously loaded programmed eBPF instructions onto the packet. At the end of the pipeline, it returns a routing decision to drop the packet, sent it to the controller, forward it to some output port, or flood it.

7.2 Hardware

7.2.1 Netronome SmartNICs. Smart network interface cards (SmartNICs) are network devices capable of having their functionality modified on runtime to implement different modes of operation. Instead of just providing connectivity and basic Layer-2 and physical layer processing, these devices can execute user-defined packet processing making use of dedicated cores and memory, freeing many cycles from the computer's CPU. To the best of our knowledge, the only company currently offering this kind of product (with support for eBPF) in the market is Netronome.

Network administrators can define packet processing routines in P4 or eBPF, for example, and load these programs to Netronome cards, which will execute the code provided upon packet reception and transmission. Different languages are supported by different firmware versions, which can be loaded to the SmartNIC without a reboot or removal of the device from the server.

With the corresponding firmware version, eBPF programs running on the TC and XDP layers can be seamlessly offloaded to these devices for increased performance. Most drivers and code necessary to perform the offloading are already part of the upstream kernel.

Given their great processing power, SmartNICs have become a good alternative for low-latency and high speed workloads. Since some functionality can be implemented directly in the hardware, packets can be modified and sent back to the network without even having to go up the operating system's network stack. Moreover, programs can be modified and updated on-the-fly. Typical use cases for such devices are early packet filtering, rate limiters, DDoS mitigation tasks, load balancing, RSS jobs, packet switching, and so on.

8 PROJECTS WITH EBPF

Although fairly recent, eBPF has already been used by many groups to power interest research- and industry-led projects. These range from tasks to support production systems operations to techniques to provide new network services. This section discusses some of these to demonstrate the wide spectrum of possible applications to eBPF.

The recent demand for greater network programmability has led to the emergence of technologies such as segment routing [29]. This technique allows network administrators to specify different actions to be performed on packets at specific points in the network. Using MPLS labels or the IPv6 protocol with a special field called Source Routing Header (SRH), each packet is encapsulated with an ordered list of routing and processing actions, called segments. As the packet is moved over the network, enabled devices process the list of segments and perform the actions specified by it. This allows the implementation, for example, of different network functions [4].

The Linux kernel already supports segment routing over IPv6 since version 4.10 but with only a few processing options, such as routing and sending and receiving labeled packets. Reference [24] used the eBPF framework to make the creation and specification of new segments more generic and flexible. Through the addition of new helper functions, the authors extended the Linux kernel to allow the implementation of segments in the form of eBPF programs. Thus, new network functions can be easily developed and associated with Linux routing rules, making it easier to integrate with segment-based routing over IPv6 environments.

Several other projects that make use of the eBPF framework to add programmability to the data plane in different contexts. InKeV [5] is a network virtualization platform that uses eBPF programs to modify the data path of virtual data center networks. To solve the OpenFlow fixed

matching problem, Reference [33] proposes to utilize eBPF programs to provide varied matching fields. Reference [64] uses eBPF to create an extensible datapath architecture to the Open vSwitch virtual switch. Reference [12] proposes a new version of the *iptables* tool using this technology.

In IoT scenarios and the edge computing paradigm, the data collected from sensors can be requested by several entities. In this case, the information is duplicated. In Reference [6], an eBPF program controls the packet duplication operation.

Some companies already use eBPF in production environments, such as Cloudflare, which uses XDP for denial-of-service attack mitigation, load balancing, and eBPF on upper layers for socket filtering and dispatching [45]. Another example comes from Facebook, which developed an L4 load balancer based on eBPF, called Katran [26]. Also, Netflix has been using eBPF for performance monitoring and system profiling [36].

Moreover, eBPF technology has been playing an important role in container networking. The veth interface type, which is commonly used on Linux for communication between containers received support for Native XDP on kernel 4.14 [41]. The Cilium [20] open-source project uses eBPF extensively to provide networking and security for microservice applications, being aware of higher-level details than just network headers. By loading eBPF programs to containers, Cilium can apply per-container security and networking policies. Weave Scope is another project that has been leveraging eBPF for a while to track TCP connections on Kubernetes clusters [69]. Project Calico has also announced earlier this year that a new data plane for container networking based on eBPF is being developed [57].

Finally, eBPF-based open source projects have also emerged in recent years. Cilium [18] is a project to provide security in container networks and microservice applications. The IOVisor [32] project maintains multiple eBPF-based subprojects such as *gobpf* [30], which allows interaction with the eBPF system using the Go language, *ply* [56] and *bpfftrace* [16] for kernel introspection, in addition to BCC [7] and *uBPF* [65], discussed previously.

9 LIMITATIONS AND WORKAROUNDS

eBPF is a powerful technology for fast packet processing and kernel programmability. However, to execute inside the kernel, some restrictions are applied to eBPF programs to guarantee system stability and security. This section discusses some eBPF limitations and workarounds to overcome them. Some of the workarounds described here can be found on the XDP project development repository available in Reference [3].

9.1 Subset of C Language Libraries

eBPF uses a restricted number of C language libraries and does not support operations with external libraries. An alternative to overcome this limitation is to define and use auxiliary functions. For example, eBPF programs cannot use the `printf` function, because they run inside the kernel and this function is not implemented in the eBPF. However, they can use the `bpf_trace_printk()` helper function, which saves log messages generated by eBPF programs according to user-defined output in the kernel trace folder (`/sys/kernel/debug/tracing/trace`). By using the log generated, the user can analyze and find possible errors in the execution of the eBPF program.

9.2 Non-static Global Variables

Currently, eBPF programs only support static global variables [13]. An alternative is to use the `BPF_MAP_TYPE_PERCPU_ARRAY` map. This map reserves a user-defined size non-shared memory space that can be used to store temporary data with a single entry during program execution [19].

9.3 Loops

983

The eBPF verifier did not allow loops to make sure all programs finish. The first technique adopted to bypass this limitation was to use loop unrolling directives from the clang compiler to rewrite a loop as a repeated sequence of independent instructions. This technique partially solves this limitation, as it can only be used when the number of repetitions can be determined at compile time. Besides that, it has a side effect of increasing the number of instructions in the final program. The code snippet below demonstrates how to tell the compiler to unroll a for loop:

```
#pragma clang loop unroll (full)
for (int i=0; i<8; i++) { ... }
```

9.3.1 Bounded Loops. In kernel version 5.3 the use of unroll loop directives was replaced by bounded loops. Since this release, the verifier does not reject all eBPF programs that contain loops without first checking if loops finish within a timeout. Bounded loops were proposed by John Fastabend and presented during the BPF Microconference at the 2018 Linux Plumbers Conference. This technique allows using simple loops modeled through the verifier. It analyzes the behavior of a loop through its induction variable and checks if memory accesses using the induction variable belong to the range of memory addresses. Implementation details about bounded loops on the verifier are available in Reference [27].

9.4 Limited Stack Space

998

Local variables of a C program are stored in the stack after its translation to an eBPF program. As the stack space is limited to only 512 bytes, it may be insufficient to store all the local variables of a program after its translation. The workaround adopted is the same used for global variables: to use the BPF_MAP_TYPE_PERCPU_ARRAY as an auxiliary buffer to store some local variables when the stack space is not enough [19].

9.5 Complex Applications

1004

Miano et al. [47] provide an in-depth discussion about the challenges faced when implementing complex network functions with eBPF. For example, applications that need to send the same packet to multiple interfaces (e.g., flood operation on a switch) are hard to implement, since the program would have to loop through all interfaces and copy the packet to each one of them. These are difficult tasks given the current available eBPF mechanisms.

In addition, since eBPF programs are associated with hooks, they follow a passive event-based model. This makes it difficult to perform asynchronous tasks, such as active network measurements, which (currently) would have to be performed by a userspace application, for example.

Moreover, the existing tools to implement a control plane for eBPF programs (e.g., libbpf and the bpf syscall) are very basic, mainly relying on map structures for data exchange. However, projects like Cilium and BCC represent progress in that respect.

Finally, the maximum number of instructions of eBPF programs used to be limited to 4096, and this diffculted the development of complex network functions. One way to get around this limitation was to split a program into multiple subprograms and jump from one subprogram to another using the bpf_tail_call() helper function. This technique enables the development of network services as a collection of loosely coupled modules, where each module performs a different function (analysis, classification, or field modification), with low overhead when jumping from one module to another. However, the maximum number of nested tail calls allowed is 32. In April 2019, the maximum number of instructions was increased from 4096 to 1 million, allowing the execution of larger eBPF programs without requiring tail calls [1].

1025 10 COMPARISON WITH SIMILAR TECHNOLOGIES

1026 As shown throughout this work, eBPF is a powerful tool to add programmability to networking
 1027 platforms and an alternative to achieve fast packet processing on Linux environments through
 1028 the XDP hook. In this section, we discuss how eBPF fits into programmable and high-speed data
 1029 planes landscapes by comparing it with some widely used technologies: P4, Click, Netmap, and
 1030 Data Plane Development Kit (DPDK).

1031 10.1 Programmable Data Planes

1032 The P4 language [14] has been proposed to enable the definition of custom data planes for pro-
 1033 grammable switches. With widespread adoption, P4-enabled devices could compose a fully pro-
 1034 grammable core network, defining the functionality of both physical and virtual switches. On the
 1035 same front, BPFabric uses eBPF as its language to define the behavior of a programmable virtual
 1036 switch, applicable to virtualized environments or even server-centric networks.

1037 However, the main strength of eBPF resides in its application to the edges of communication.
 1038 In the Linux kernel, it is capable of defining data plane functionalities at the communication end-
 1039 points, a different use case than the one P4 is designed for. Switching and routing are done by the
 1040 core network, covered by P4, but a significant portion of the network protocol stack is implemented
 1041 at the endpoints. With eBPF, this remaining part can be monitored, modified and reconfigured on
 1042 demand, making it an essential tool to provide full network programmability.

1043 Click [37] also enables the creation of custom network elements to process packets inside the
 1044 Linux kernel. Applications can be created through the composition of many components, called
 1045 elements, which could also be implemented as C++ code using Click-specific function calls. The
 1046 definition of a Click application is done through a configuration file specifying the kinds of ele-
 1047 ments used and the interconnection between them, representing a packet processing graph. This
 1048 specification can then be compiled to userspace or kernel. Click kernel modules capture packets
 1049 close to the network device, and can pass them to the kernel stack using *ToHost* elements, similar to
 1050 what can be done in the XDP hook with eBPF. Li et al. [38] build upon Click to create the ClickNP
 1051 framework that enables the creation and execution of network functions on FPGAs, achieving
 1052 40 Gbps throughput for any packet size.

1053 Although Click and ClickNP provide a broader set of pre-built primitives to create routers and
 1054 network functions, they do not offer the same level of integration with the kernel stack as eBPF.
 1055 The latter allows attaching programs to interact with several layers of the kernel stack, providing a
 1056 higher degree of control over the kernel's packet processing mechanisms. However, Click could be
 1057 modified to take advantage of the eBPF/XDP infrastructure to be its basis to interact with packet
 1058 processing facilities, and possibly combine its expressiveness with eBPF's performance and native
 1059 kernel integration.

1060 10.2 High-speed Packet Processing

1061 On high-speed networks with 10 Gbps links and beyond, inter-packet arrival times can get as
 1062 low as tens of nanoseconds, leaving very little time to process each packet. Due to this stringent
 1063 requirement, common system operations become too costly, such as context switches and inter-
 1064 rupt handling. Well-known technologies like DPDK [42] and NetMap [58] handle this problem by
 1065 operating in poll mode and bypassing the kernel altogether, performing all packet processing in
 1066 user space. By the use of specialized drivers, packets received by the NIC are sent directly to a
 1067 user space application, which will process them. In addition, the DPDK library and applications
 1068 built with it are usually optimized for aligned memory access, multi-core processing, non-uniform
 1069 memory access (NUMA) and other optimizations aimed to save precious CPU cycles.

Fast Packet Processing with eBPF and XDP

16:33

The eBPF system in the kernel provides good performance in a different manner: by allowing custom packet processing at the XDP hook, which is the lowest kernel layer. Through this hook, eBPF applications can parse, modify, collect statistics, and take action on incoming packets possibly without going through the OS' network stack, forwarding packets directly back to the network. This way, context switches are avoided by embedding all network processing in the kernel.

If compared to DPDK, then one of the advantages of the XDP hook is being able to use existing network facilities present in the kernel (ex: routing tables), which have to be re-implemented from the ground up when processing is done in user space, which is the case of DPDK. Also, since XDP is handled by the kernel, it can benefit from kernel API stability guarantees, from the separation mechanisms already in place to enforce security, integration with the existing stack without requiring re-injecting packets through an exception path. Device sharing is also facilitated, as the program does not need full-control over the NIC, which can remain visible to the OS. At the same time, programs are transparent to other applications on the host, since processing is hidden under the OS' abstraction layers. In terms of resource usage, it avoids investing precious CPU cycles with busy-polling due to its event-based nature. Recent performance comparisons between the two technologies show that DPDK still reaches higher bandwidths (115 Mpps against 100 Mpps for XDP when dropping packets with five cores), but at the cost of a much higher CPU usage than XDP [31]. Finally, programs can be replaced atomically, providing greater flexibility for on-demand changes to packet processing.

Thus, eBPF and DPDK provide different approaches to high-speed packet processing, both with their own set of restrictions and advantages, as discussed in this section and in Section 9. The best choice of technology may depend on the actual use case, although integration of both could be achieved by using AF_XDP sockets and the corresponding poll mode driver on DPDK [22] or the *librte_bpf* library provided by DPDK to add an eBPF VM directly to DPDK applications [23].

11 CONCLUSION

In this work, we presented a vision of the theoretical and practical aspects related to fast packet processing with eBPF and XDP. In the theoretical part, we discussed the BPF and eBPF machines, an overview of the eBPF system provided by the Linux kernel, the available hooks and some results of recent research. In the practical part, we focused on eBPF and the XDP hook, providing examples and showing existing tools.

Given their power for fast packet processing, we consider that there is great potential in the development of new research projects with eBPF and XDP, either as a tool for the development of new network functions, or allowing to provide new functionalities in the data plane, such as the creation of new communication standards and protocols, or in the development of new research prototypes and network solutions. Together, eBPF and XDP will help develop new interesting research with great potential in the area of Computer Networking.

REFERENCES

- [1] 2019. bpf: Increase Complexity Limit and Maximum Program Size. Retrieved from <https://git.kernel.org/pub/scm/linux/kernel/git/davem/net-next.git/commit/?id=c04c0d2b968ac45d6ef020316808ef6c82325a82>.
- [2] 2019. libbpf Unification and Golang Bindings. *Discussion Summary, Linux Kernel Developers' bpfconf 2019*. Retrieved from <http://vger.kernel.org/bpfconf2019.html#session-4>.
- [3] 2019. XDP Project Repository. Retrieved from <https://github.com/xdp-project/xdp-project>.
- [4] Ahmed Abdelsalam, Francois Clad, Clarence Filsfils, Stefano Salsano, Giuseppe Siracusano, and Luca Veltri. 2017. Implementation of virtual network function chaining through segment routing in a linux-based NFV infrastructure. In *Proceedings of the 2017 IEEE Conference on Network Softwarization: Softwarization Sustaining a Hyper-Connected World: en Route to 5G (NetSoft'17)*. IEEE, Los Alamitos, CA, 1–5. DOI: <https://doi.org/10.1109/NETSOFT.2017.8004208> arXiv:1702.05157

Q1

- 1116 [5] Zaafar Ahmed, Muhammad Hamad Alizai, and Affan A. Syed. 2018. InKeV: In-kernel distributed network virtual-
 1117 ization for DCN. *SIGCOMM Comput. Commun. Rev.* 46, 3, Article 4 (Jul. 2018), 6 pages. DOI : [https://doi.org/10.1145/](https://doi.org/10.1145/3243157.3243161)
 1118 [3243157.3243161](https://doi.org/10.1145/3243157.3243161)
- Q2 1119 [6] S. Baidya, Y. Chen, and M. Levorato. 2018. eBPF-based content and computation-aware communication for real-time
 1120 edge computing. In *Proceedings of the INFOCOM IEEE Conference on Computer Communications Workshops (INFOCOM*
 1121 *WORKSHOPS'18)*. IEEE, Los Alamitos, CA, 865–870. DOI : <https://doi.org/10.1109/INFCOMW.2018.8407006>
- 1122 [7] BCC. 2019. BPF Compiler Collection. Retrieved from <https://github.com/iovisor/bcc>.
- 1123 [8] BCC. 2019. BPF Program Types. Retrieved from [https://github.com/iovisor/bcc/blob/master/docs/kernel-versions.](https://github.com/iovisor/bcc/blob/master/docs/kernel-versions.md#program-types)
 1124 [md#program-types](https://github.com/iovisor/bcc/blob/master/docs/kernel-versions.md#program-types).
- 1125 [9] BCC. 2019. XDP Compatible Drivers. Retrieved from [https://github.com/iovisor/bcc/blob/master/docs/](https://github.com/iovisor/bcc/blob/master/docs/kernel-versions.md#xdp)
 1126 [kernel-versions.md#xdp](https://github.com/iovisor/bcc/blob/master/docs/kernel-versions.md#xdp).
- Q3 1127 [10] David Beckett, Jaco Joubert, and Simon Horman. 2018. Host Dataplane Acceleration (HDA).
- 1128 [11] Gilberto Bertin. 2017. XDP in practice: Integrating XDP into our DDoS mitigation pipeline. In *Proceedings of the*
 1129 *Netdev 2.1 Technical Conference on Linux Networking*. 1–5.
- 1130 [12] Matteo Bertrone, Sebastiano Miano, Fulvio Rizzo, and Massimo Tumolo. 2018. Accelerating linux security with eBPF
 1131 iptables. In *Proceedings of the ACM SIGCOMM 2018 Conference on Posters and Demos (SIGCOMM'18)*. ACM, New York,
 1132 NY, 108–110. DOI : <https://doi.org/10.1145/3234200.3234228>
- 1133 [13] Daniel Borkmann. 2019. bpf, Libbpf: Support Global Data/bss/rodata Sections. Retrieved from [https://git.kernel.org/](https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/commit/?id=d859900c4c56)
 1134 [pub/scm/linux/kernel/git/torvalds/linux.git/commit/?id=d859900c4c56](https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/commit/?id=d859900c4c56).
- 1135 [14] Pat Bosshart, Dan Daly, Glen Gibb, Martin Izzard, Nick McKeown, Jennifer Rexford, Cole Schlesinger, Dan Talayco,
 1136 Amin Vahdat, George Varghese, and David Walker. 2014. P4: Programming protocol-independent packet processors.
 1137 *SIGCOMM Comput. Commun. Rev.* 44, 3 (Jul. 2014), 87–95. DOI : <https://doi.org/10.1145/2656877.2656890>
- 1138 [15] Autores bpftool. 2018. Manual Bpftool. Retrieved from [https://elixir.bootlin.com/linux/v4.18-rc1/](https://elixir.bootlin.com/linux/v4.18-rc1/source/tools/bpf/bpftool/Documentation/bpftool.rst) [source/tools/bpf/](https://elixir.bootlin.com/linux/v4.18-rc1/source/tools/bpf/bpftool/Documentation/bpftool.rst)
 1139 [bpftool/Documentation/bpftool.rst](https://elixir.bootlin.com/linux/v4.18-rc1/source/tools/bpf/bpftool/Documentation/bpftool.rst).
- 1140 [16] bpfftrace. 2019. High-level Tracing Language for Linux eBPF. Retrieved from <https://github.com/iovisor/bpfftrace>.
- 1141 [17] Mihai Budiu. 2015. Compiling p4 to ebpf. Retrieved from [https://github.com/iovisor/bcc/tree/master/src/cc/frontends/](https://github.com/iovisor/bcc/tree/master/src/cc/frontends/p4)
 1142 [p4](https://github.com/iovisor/bcc/tree/master/src/cc/frontends/p4).
- 1143 [18] Cilium. 2018. Cilium 1.0: Bringing the BPF Revolution to Kubernetes Networking and Security. Retrieved from
 1144 <https://cilium.io/blog/2018/04/24/cilium-10/>.
- 1145 [19] Cilium. 2019. BPF and XDP Reference Guide. Retrieved September 9, 2019 from [https://cilium.readthedocs.io/en/](https://cilium.readthedocs.io/en/latest/bpf/)
 1146 [latest/bpf/](https://cilium.readthedocs.io/en/latest/bpf/).
- 1147 [20] Cilium. 2019. Cilium: API-aware Networking and Security. Retrieved September 9, 2019 from <https://cilium.io/>.
- 1148 [21] Jonathan Corbet. 2014. BPF: The Universal In-kernel Virtual Machine. Retrieved from [https://lwn.net/Articles/](https://lwn.net/Articles/599755/)
 1149 [599755/](https://lwn.net/Articles/599755/).
- 1150 [22] DPDK. 2019. AF_XDP Poll Mode Driver. Retrieved from https://doc.dpdk.org/guides/nics/af_xdp.html.
- 1151 [23] DPDK. 2019. Berkeley Packet Filter Library. Retrieved from https://doc.dpdk.org/guides/prog_guide/bpf_lib.html.
- 1152 [24] Fabien Duchene, Mathieu Jadin, and Olivier Bonaventure. 2018. Exploring various use cases for IPv6 segment routing.
 1153 In *Proceedings of the ACM SIGCOMM 2018 Conference on Posters and Demos (SIGCOMM'18)*. ACM, New York, NY, 129–
 1154 131. DOI : <https://doi.org/10.1145/3234200.3234213>
- 1155 [25] Eric Dumazet. 2011. A JIT for Packet Filters. Retrieved from <https://lwn.net/Articles/437981/>.
- 1156 [26] Facebook. 2018. Katran Source Code Repository. Retrieved from <https://github.com/facebookincubator/katran>.
- 1157 [27] John Fastabend. 2018. [RFC PATCH 00/16] bpf, Bounded Loop Support Work in Progress. Retrieved from [https://lwn.](https://lwn.net/ml/netdev/20180601092646.15353.28269.stgit@john-Precision-Tower-5810/)
 1158 [net/ml/netdev/20180601092646.15353.28269.stgit@john-Precision-Tower-5810/](https://lwn.net/ml/netdev/20180601092646.15353.28269.stgit@john-Precision-Tower-5810/).
- 1159 [28] Nick Feamster, Jennifer Rexford, and Ellen Zegura. 2014. The road to SDN: An intellectual history of programmable
 1160 networks. *ACM SIGCOMM Comput. Commun. Rev.* 44, 2 (2014), 87–98.
- 1161 [29] C. Filsfils, S. Previdi, L. Ginsberg, B. Decraene, S. Litkowski, and R. Shakir. 2018. *Segment Routing Architecture*. RFC
 1162 8402. RFC Editor.
- 1163 [30] gobpf. 2019. Go Bindings for Creating BPF Programs. Retrieved from <https://github.com/iovisor/gobpf>.
- 1164 [31] Toke Høiland-Jørgensen, Jesper Dangaard Brouer, Daniel Borkmann, John Fastabend, Tom Herbert, David Ahern, and
 1165 David Miller. 2018. The eXpress data path: Fast programmable packet processing in the operating system kernel. In
 1166 *Proceedings of the 14th International Conference on Emerging Networking EXperiments and Technologies (CoNEXT'18)*.
 1167 ACM, New York, NY, 54–66. DOI : <https://doi.org/10.1145/3281411.3281443>
- 1168 [32] IOvisor. 2019. Iovisor Project. Retrieved March 29, 2019 from www.iovisor.org.
- 1169 [33] S. Jouet, R. Cziva, and D. P. Pezaros. 2015. Arbitrary packet matching in OpenFlow. In *Proceedings of the 16th Inter-*
 1170 *national Conference on High Performance Switching and Routing (HPSR'15)*. IEEE, Los Alamitos, CA, 1–6. DOI : [https://](https://doi.org/10.1109/HPSR.2015.7483106)
 1171 doi.org/10.1109/HPSR.2015.7483106

Fast Packet Processing with eBPF and XDP

16:35

- [34] Simon Jouet and Dimitrios P. Pezaros. 2017. BPFabric: Data plane programmability for software defined networks. 1172
In *Proceedings of the Symposium on Architectures for Networking and Communications Systems (ANCS'17)*. IEEE Press, 1173
Piscataway, NJ, 38–48. DOI: <https://doi.org/10.1109/ANCS.2017.14> 1174
- [35] Michael Kerrisk. 2013. BFP 8 Linux Manual Page. Retrieved June 8, 2019 from <http://man7.org/linux/man-pages/man8/bpf.8.html>. 1175
- [36] Jason Koch, Martin Spier, Brendan Gregg, and Ed Hunter. 2019. Extending Vector with eBPF to Inspect Host and Con- 1177
tainer Performance. Retrieved from [https://medium.com/netflix-techblog/extending-vector-with-ebpf-to-inspect-](https://medium.com/netflix-techblog/extending-vector-with-ebpf-to-inspect-host-and-container-performance-5da3af4c584b) 1178
[host-and-container-performance-5da3af4c584b](https://medium.com/netflix-techblog/extending-vector-with-ebpf-to-inspect-host-and-container-performance-5da3af4c584b). 1179
- [37] Eddie Kohler, Robert Morris, Benjie Chen, John Jannotti, and M. Frans Kaashoek. 2000. The click modular router. 1180
ACM Trans. Comput. Syst. 18, 3 (2000), 263–297. 1181
- [38] Bojie Li, Kun Tan, Layong (Larry) Luo, Yanqing Peng, Renqian Luo, Ningyi Xu, Yongqiang Xiong, Peng Cheng, and 1182
Enhong Chen. 2016. ClickNP: Highly flexible and high performance network processing with reconfigurable hard- 1183
ware. In *Proceedings of the 2016 ACM SIGCOMM Conference (SIGCOMM'16)*. ACM, New York, NY, 1–14. DOI: <https://doi.org/10.1145/2934872.2934897> 1184
- [39] libbpf. 2018. Libbpf Source Code. Retrieved from <https://elixir.bootlin.com/linux/v4.18-rc1/source/tools/lib/bpf>. 1185
- [40] libbpf. 2019. Stand-alone Libbpf. Retrieved from <https://github.com/libbpf/libbpf>. 1186
- [41] Linux. 2017. Net: Xdp: Support Xdp Generic on Virtual Devices. Retrieved from [https://git.kernel.org/pub/scm/linux/](https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/commit/?id=d445516966dcb2924741b13b27738b54df2af01a) 1188
[kernel/git/torvalds/linux.git/commit/?id=d445516966dcb2924741b13b27738b54df2af01a](https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/commit/?id=d445516966dcb2924741b13b27738b54df2af01a). 1189
- [42] Linux Foundation. 2015. Data Plane Development Kit (DPDK). Retrieved from <http://www.dpdk.org>. 1190
- [43] D. F. Macedo, D. Guedes, L. F. M. Vieira, M. A. M. Vieira, and M. Nogueira. 2015. Programmable networks: From 1191
software-defined radio to software-defined networking. *IEEE Commun. Surv. Tutor.* 17, 2 (2015), 1102–1125. DOI: 1192
<https://doi.org/10.1109/COMST.2015.2402617> 1193
- [44] Alan Maguire. 2019. Notes on BPF (1)—A Tour of Program Types. Retrieved from [https://blogs.oracle.com/linux/](https://blogs.oracle.com/linux/notes-on-bpf-1) 1194
[notes-on-bpf-1](https://blogs.oracle.com/linux/notes-on-bpf-1). 1195
- [45] Marek Majkowski. 2019. Cloudflare Architecture and How BPF Eats the World. Retrieved from [https://blog.cloudflare.](https://blog.cloudflare.com/cloudflare-architecture-and-how-bpf-eats-the-world/) 1196
[com/cloudflare-architecture-and-how-bpf-eats-the-world/](https://blog.cloudflare.com/cloudflare-architecture-and-how-bpf-eats-the-world/). 1197
- [46] Steven McCanne and Van Jacobson. 1993. The BSD packet filter: A new architecture for user-level packet capture. 1198
In *Proceedings of the USENIX Winter 1993 Conference Proceedings on USENIX Winter 1993 Conference Proceedings* 1199
(USENIX'93). USENIX Association, Berkeley, CA, 1–11. 1200
- [47] Sebastiano Miano, Matteo Bertrone, Fulvio Risso, Massimo Tumolo, and Mauricio Vásquez Bernal. 2018. Creating 1201
complex network service with ebpf: Experience and lessons learned. In *Proceedings of the High Performance Switching* 1202
and Routing (HPSR'18). IEEE, Los Alamitos, CA, 1–8. 1203
- [48] Rashid Mijumbi, Joan Serrat, Juan Luis Gorricho, Niels Bouten, Filip De Turck, and Raouf Boutaba. 2016. Network 1204
function virtualization: State-of-the-art and research challenges. *IEEE Communi. Surv. Tutor.* 18, 1 (2016), 236–262. 1205
- [49] David Miller. 2017. BPF Verifier Overview. Retrieved April 9, 2019 from <https://lwn.net/Articles/794934/>. 1206
- [50] J. Mogul, R. Rashid, and M. Accetta. 1987. The packer filter: An efficient mechanism for user-level network code. 1207
In *Proceedings of the 11th ACM Symposium on Operating Systems Principles (SOSP'87)*. ACM, New York, NY, 39–51. 1208
DOI: <https://doi.org/10.1145/41457.37505> 1209
- [51] Quentin Monnet. 2019. All-Out Programmability in Linux: An Introduction to BPF as a Monitoring Tool. Retrieved 1210
April 9, 2019 from https://qmo.fr/docs/talk_20190516_allout_programmability_bpf.pdf. 1211
- [52] Quentin Monnet. 2019. Rust Virtual Machine and JIT Compiler for eBPF Programs. Retrieved from [https://github.](https://github.com/qmonnet/rbpf) 1212
[com/qmonnet/rbpf](https://github.com/qmonnet/rbpf). 1213
- [53] Netronome. 2019. Sample BPF Offload Apps. Retrieved from <https://github.com/Netronome/bpf-samples>. 1214
- [54] PLUMgrid. 2016. Linux Kernel Source Code. Retrieved June 7, 2019 from [https://github.com/torvalds/linux/blob/v5.](https://github.com/torvalds/linux/blob/v5.3/samples/bpf/xdp1_kern.c) 1215
[3/samples/bpf/xdp1_kern.c](https://github.com/torvalds/linux/blob/v5.3/samples/bpf/xdp1_kern.c). 1216
- [55] PLUMgrid. 2016. Linux Kernel Source Code. Retrieved June 7, 2019 from [https://github.com/torvalds/linux/blob/v5.](https://github.com/torvalds/linux/blob/v5.3/samples/bpf/xdp1_user.c) 1217
[3/samples/bpf/xdp1_user.c](https://github.com/torvalds/linux/blob/v5.3/samples/bpf/xdp1_user.c). 1218
- [56] ply. 2019. Dynamic Tracing in Linux. Retrieved from <https://github.com/iovisor/ply>. 1219
- [57] Alex Pollitt. 2019. Tigera adds eBPF Support to Calico. Retrieved September 9, 2019 from [https://www.projectcalico.](https://www.projectcalico.org/tigera-adds-ebpf-support-to-calico/) 1220
[org/tigera-adds-ebpf-support-to-calico/](https://www.projectcalico.org/tigera-adds-ebpf-support-to-calico/). 1221
- [58] Luigi Rizzo. 2012. netmap: A novel framework for fast packet I/O. In *Proceedings of the 2012 USENIX Annual Technical* 1222
Conference (USENIX ATC'12). USENIX Association, Berkeley, CA, 101–112. [https://www.usenix.org/conference/atc12/](https://www.usenix.org/conference/atc12/technical-sessions/presentation/rizzo) 1223
[technical-sessions/presentation/rizzo](https://www.usenix.org/conference/atc12/technical-sessions/presentation/rizzo) 1224
- [59] Marta Rybczyńska. 2019. Bounded Loops in BPF for the 5.3 Kernel. Retrieved April 9, 2019 from [https://www.spinics.](https://www.spinics.net/lists/xdp-newbies/msg00185.html) 1225
[net/lists/xdp-newbies/msg00185.html](https://www.spinics.net/lists/xdp-newbies/msg00185.html). 1226
- [60] Jay Schulist, Daniel Borkmann, and Alexei Starovoitov. 2019. Linux Socket Filtering aka Berkeley Packet Filter (BPF). 1227
Retrieved March 17, 2019 from www.kernel.org/doc/Documentation/networking/filter.txt. 1228

- 1229 [61] Haoyu Song. 2013. Protocol-oblivious forwarding: Unleash the power of SDN through a future-proof forward-
 1230 ing plane. In *Proceedings of the Second ACM SIGCOMM Workshop on Hot Topics in Software Defined Networking*
 1231 (*HotSDN'13*). ACM, New York, NY, 127–132. DOI : <https://doi.org/10.1145/2491185.2491190>
- 1232 [62] Alexei Starovoitov. 2015. bpf: Introduce bpf_tail_call() Helper. Retrieved from <https://lwn.net/Articles/645169/>.
- 1233 [63] Alexei Starovoitov. 2018. Lifetime of BPF Objects. Retrieved from [https://facebookmicrosites.github.io/bpf/blog/2018/](https://facebookmicrosites.github.io/bpf/blog/2018/08/31/object-lifetime.html)
 1234 [08/31/object-lifetime.html](https://facebookmicrosites.github.io/bpf/blog/2018/08/31/object-lifetime.html).
- 1235 [64] Cheng-Chun Tu, Joe Stringer, and Justin Pettit. 2017. Building an extensible open vSwitch datapath. *SIGOPS Oper.*
 1236 *Syst. Rev.* 51, 1 (Sep. 2017), 72–77. DOI : <https://doi.org/10.1145/3139645.3139657>
- 1237 [65] uBPF. 2019. Userspace eBPF VM. Retrieved from <https://github.com/iovisor/ubpf>.
- 1238 [66] Marcos A. M. Vieira, Matheus S. Castanho, Racyus D. G. Pacífico, Elerson R. S. Santos, Eduardo P. M. Câmara
 1239 Júnior, and Luiz F. M. Vieira. 2019. Zenodo—eBPF Tutorial. Retrieved from [https://zenodo.org/record/3519347#.](https://zenodo.org/record/3519347#.XbMxR6zMNhE)
 1240 [XbMxR6zMNhE](https://zenodo.org/record/3519347#.XbMxR6zMNhE).
- 1241 [67] Marcos A. M. Vieira, Matheus S. Castanho, Racyus D. G. Pacífico, Elerson R. S. Santos, Eduardo P. M. Câmara Júnior,
 1242 and Luiz F. M. Vieira. 2019. eBPF Tutorial. Retrieved from <https://github.com/racyusdelano0/bpf-tutorial>.
- 1243 [68] VMWare. 2018. p4c-xdp. Retrieved from <https://github.com/vmware/p4c-xdp>.
- 1244 [69] WeaveWorks. 2017. Improving Performance and Reliability in Weave Scope with eBPF. Retrieved from [https://www.](https://www.weave.works/blog/improving-performance-reliability-weave-scope-ebpf/)
 1245 [weave.works/blog/improving-performance-reliability-weave-scope-ebpf/](https://www.weave.works/blog/improving-performance-reliability-weave-scope-ebpf/).
- 1246 [70] XDP-Project. 2019. AXDP Hands-On Tutorial. Retrieved August 20, 2019 from [https://github.com/xdp-project/](https://github.com/xdp-project/xdp-tutorial)
 1247 [xdp-tutorial](https://github.com/xdp-project/xdp-tutorial).

1248 Received June 2019; revised October 2019; accepted October 2019

Author Queries

Q1: AU: Refs 1-3: Please provide authors for all refs.

Q2: AU: Please verify Retrieved from dates; the style changes throughout the list.

Q3: AU: Ref. 10: Please update and complete.