

Server agnostic DNS augmentation using eBPF

Master thesis by Tom Carpay
Security and Network Engineering masters programme

Supervisors: Willem Toorop and Luuk Hendriks
NLnet Labs

August 17, 2020

With ever larger DNS request volumes, handling requests efficiently becomes ever more important. One way to address higher network performance is to bypass the kernel network stack completely, however the software providing the service then has the daunting task to perform all the low-level tasks the kernel would normally handle. The Extended Berkeley Packet Filter (eBPF) and in particular Express Data Path (XDP) kernel hook provides an alternative approach, in which traffic can be considered and acted upon, or passed on to the kernel, from the lowest layer of the network stack.

This research seeks to take this approach one step further and augment DNS software agnostic of the software providing the (basic) DNS service. To this end we explore several aspects of the XDP kernel hook introduced by eBPF. We examine two use cases: QName rewriting and Response Rate Limiting (RRL). We develop a prototype for QName rewriting and develop three RRL prototypes, including an augmentation that is currently unavailable in other DNS software. Each of these prototypes demonstrates one or more different XDP eBPF functionalities to augment DNS. The QName rewriting prototype is used in a continued setup involving the RIPE Atlas network. We perform experiments for each Response Rate Limit prototype by stress testing an authoritative DNS server and introducing the prototypes to the setup. We show that the CPU load drops when using the prototypes over their user space counterparts and we conclude that XDP eBPF is a viable candidate for augmentations to DNS software agnostic of the supplier.

Contents

1	Introduction	3
2	Related work	5
3	Background	5
4	Methods	8
4.1	QName rewriting	9
4.2	Response Rate Limiting	9
5	Experiments	10
5.1	QName rewrite	11
5.2	General RRL	11
5.3	Per IP RRL	12
5.4	Unknown host RRL	13
6	Results	13
6.1	QName rewrite	13
6.2	General RRL	14
6.3	Per IP RRL	16
6.4	Unknown host RRL	17
7	Discussion	18
8	Conclusion	19
9	Future work	20
10	Acknowledgements	22
A	QName rewrite prototype code	23
B	General RLL prototype code	36
C	Per IP RRL prototype code	43
D	Unknown host RRL prototype code	51

1 Introduction

As the internet is constantly growing, the amount of DNS requests grows with it. With ever larger request volumes, handling requests more efficiently becomes more and more of a necessity. This increase in volume becomes apparent in the packet processing limitations of network stacks at packet rates higher than 10 Gbps [1].

Currently, there is no Domain Name System (DNS) interaction available low in the Linux kernel network stack. This is desirable for operators of DNS software, for example at high volume authoritative nameservers. A method for handling DNS messages more efficiently is by removing the overhead from the network stack, which could significantly improve performance. A candidate technology this research explores for this is the Berkely Packet Filter (BPF).

The BPF is a Linux based, in-kernel virtual machine that runs a custom 64 bit instruction set that is Just-In-Time (JIT) compiled to machine instructions. BPF allows compiled C BPF bytecode to be attached to a data path via the built in verifier. Once the program passes the verifying process, which is discussed in more detail in Section 3, the program is executed whenever the attached path is traversed.

The extended Berkely Packet Filter (eBPF), the 2014 expansion added to the retroactively called classic Berkely Packet Filter (cBPF), offers a significant number of improvements which enhance its capabilities as packet filter and processor, and shift the focus to a low level tracing tool set. One of the aspects that enables this research is the addition of the eXpress Data Path (XDP). This high performance kernel hook allows BPF programs to be attached to a network interface and execute instructions on a per packet basis. Instructions can be hardware offloaded on selected Network Interface Cards (NICs). Further improvements are discussed in more detail in Section 3. BPF is native to the Linux kernel, so BPF programs are widely available. From Section 2 we learn that BPF programs are more performant to their user space counterparts in many cases.

To explore the capabilities of this technology in relation to DNS, we create a number of prototypes with limited functionalities. These functionalities are chosen based on operator needs and functionalities that are not currently available in DNS software, or are improved by this technology. If successful, these prototypes implement their respective functionalities agnostic of the software supplier, as the logic that is executed is completely invisible to user space. If successful, further developed and be used to augment all current DNS software.

We examine two main use cases: a Query Name (QName) rewriting experimental setup, and Response Rate Limiting, the DNS addition to mitigate DNS amplification attacks. For both use cases we explore prototypes to examine the functionalities of BPF.

The QName rewrite prototype is used in an experimental setup for research being done by Koolhaas and Slokker [2]. For their measurements, a method

is needed to change the DNS query to invoke a predetermined sized response. The authoritative DNS server on the test setup replies different sized responses based on the query name. The probes from which the queries are sent, embed their (variable sized) probe IDs in the query as these are needed to process the measurement results. Besides rewriting the query name to request different sized responses, also the length of the probe ID within the query name needs to be compensated, to guarantee the same (requested) sized responses for all probes. As we show in the Section 4, BPF offers a viable solution to this need, as the measurements can be taken without the need for a DNS parser or modification to a DNS software, facilitating easier reproducibility of the experiment.

The other prototypes revolve around Response Rate Limiting (RRL), a DNS enhancement to mitigate DNS amplification attacks [3]. While RRL is a functionality that is present in a number of the well known DNS software products (BIND¹, NSD², KnotDNS³) it is not present in others. Software native RRL implementations leave the overhead of the network stack in place for every packet.

The RRL prototypes examine the possibilities for RRL using BPF in three iterations. The first prototype implements a general version of RRL for all incoming DNS messages. The second prototype differentiates the RRL “buckets” for individual IP addresses. The final prototype implements a novel method where a configurable allowlist of “known senders” are exempt from RRL, while all other senders are subject to it per individual source IP address. With each prototype we explore the different functionalities that BPF offers for DNS.

The prototypes for both uses cases are experimentally verified to be functional. In case of the QName rewriting, the XDP program is verified to work on a small scale setup before being used in the large scale experiment by Koolhaas and Slokker.

With the different RRL prototypes, we examine the impact on performance that the XDP programs have to their user space counterpart, where possible. This comparison is drawn by observing the change in CPU load when both cases are being subjected to a stress test. If a comparison is not possible, as with the final RRL prototype iteration, we examine how the prototype scales within the limitation of the resources of this research.

The aim of this research is to augment DNS software agnostic of the software supplier. To this end we pose the main research question: How can XDP eBPF be used to augment and improve DNS software?

To help answer the main research question, we pose the following research sub-questions:

- Which features from XDP BPF could be used to augment DNS software?
- How can DNS augmentations be implemented based upon these XDP

¹Berkely Internet Name Domain: <https://www.isc.org/bind/>

²NLnet Labs’ open source Name Server Daemon: <https://www.nlnetlabs.nl/projects/nsd/about/>

³Open source authoritative-only name server: <https://www.knot-dns.cz>

eBPF features?

- How do these implementations impact performance?

This paper is structured as follows. In Section 2 we discuss the related work of this research and in Section 3 we discuss the concepts of BPF and XDP in more detail. In Section 4 we examine the research methods, prototypes and algorithms, and in Section 5 the experiments. In Section 6 we examine the results of the experiments and in Section 7 we discuss our findings. Finally, in Section 8 we show our conclusions and in Section 9 we discuss the future work for this research.

2 Related work

The eBPF instruction set itself has been a thoroughly researched topic. In 2017 Vieira et al. [4] gave comprehensive insight into eBPF and its data path hooks. In 2020, Høiland-Jørgenson et al. [5] extended this insight by focusing on the XDP kernel hook. They also compare the packet drop performance of XDP compared to its user space counterparts and show that XDP is more efficient. Their research builds a fundamental understanding of eBPF and the XDP kernel hook, and their conclusion on XDP efficiency build the foundations for this research.

A networking use case for XDP eBPF is for Distributed Denial Of Service (DDOS) protection. This has been studied in scientific research [6] [7] as well as by large organisations such as Cloudflare. The research has shown that an XDP eBPF DDOS program was possible technically and functional [8] and show their product works well [9]. The findings from both the industry and scientific sources show that an XDP protection program offers more efficient DDOS protection than its user space counterparts. Our Response Rate Limiting prototype aims to extend the use case to DNS, to provide protection against the DNS Amplification Reflection type of Denial of Service attacks, independent of the DNS software in use in user space.

3 Background

Originally BPF was designed as a low level packet filter with a customisable set of rules. With the follow up of eBPF in 2014 the functionality of the capabilities of the technology have been expanded to include network performance, firewalls, security, tracing, and device drivers [10].

For clarity, we distinguish classic BPF and extended BPF, which are retroactively named cBPF and BPF respectively, as in the rest of this research we adhere to this naming convention.

The extension on the cBPF improves on several points as shown in Table 1. One of the changes is that the extension of BPF introduces the ability of executing a limited set of kernel functions through “helper functions”, which was not allowed in cBPF. An example of a helper function is `bpf_ktime_get_ns`, which calls the kernel function `ktime_get_ns`. Only kernel functions that have a

Property	cBPF	eBPF
Number of registers	2	11
Register size	32 bits	64 bits
Stack size	16 bytes	512 bytes
Instruction execution limit	10.000	1.000.000
Kernel function calls	Not allowed	Allowed
Program offloading to SmartNIC	Not possible	Possible
Dynamic loading and program reloading	Not possible	Possible

Table 1: The changes made extending BPF [4].

helper function counterpart can be used in BPF programs. Specifically, an BPF program cannot call user-space functions, because it is not running in user space.

BPF programs are user-created programs that run in kernel space. Code written in C is compiled to bytecode as a BPF program. Running user-created programs in kernel space could potentially lead to security or stability risks. To mitigate this risk, all BPF programs are inspected by the BPF verifier. If the program passes the verifier, on first execution the bytecode program is JIT compiled to machine instructions to be executed. This life cycle is shown in Figure 1.

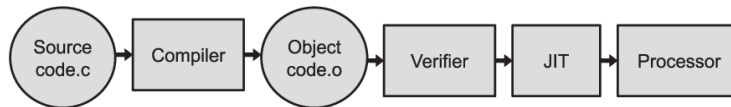


Figure 1: The life cycle of a BPF Program (copied from Vieira et al. [4]).

The BPF verifier performs three main checks [11]. The first check tries to verify that the program does not contain any kernel locking loops and that the program terminates. To this end, the verifier creates a directed acyclic Control Flow Graph (CFG) and does a depth first search of this graph to find the terminating points. If it finds any potentially non terminating path, the verifier will fail the program, as these are not allowed. An example CFG is shown in figure 2.

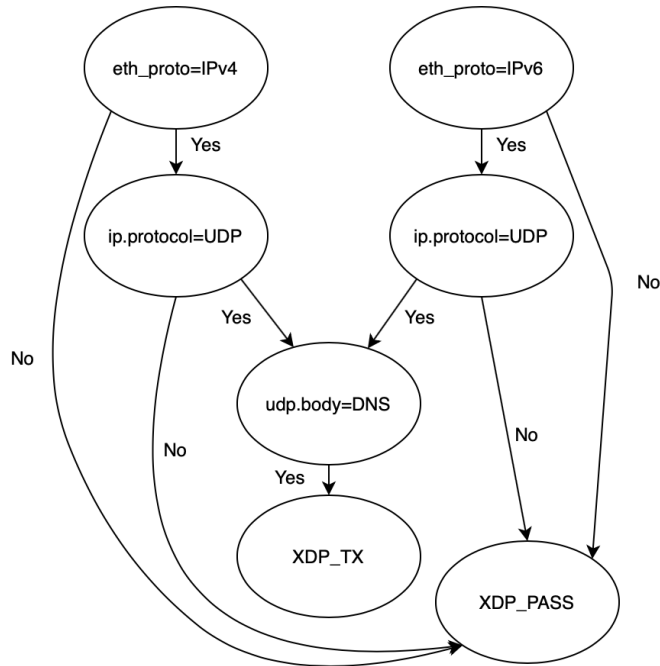


Figure 2: A directed acyclic control flow graph representation for a BPF Program that rejects all incoming DNS queries. This program follows a different flow based on the received IP version. Note the distinction between the starting points for both IP versions. This distinction is made to ensure that all instructions are reachable in either IP case during verification.

The second verifier check involves simulating the program, executing and checking per single instruction. This simulation verifies that the program stays within bounds and never accesses out of range memory.

For the final verifier check, the verifier restricts kernel functions and data structures available to the program according to the eBPF program type. In this research we use program types corresponding to the XDP kernel hook. The XDP hook is limited to incoming traffic. A more elaborate description of the verifier is described by Vieira et al. [4] and Høiland-Jørgensen et al. [5].

The method that BPF uses to interact with user space is via BPF “maps”. These maps are data structures of which the type has to be defined at compile time and are stored in user space. The map types are used in the final verifier check to restrict map usage to its respective type. BPF programs can interact with maps through BPF helper functions such as `bpf_map_lookup_elem()`, `bpf_map_update_elem()`, and `bpf_map_delete_elem()`. Since BPF programs attach to a data path which can be traversed many times per second, to avoid memory locks the map types can be defined as per CPU core or shared between all CPU cores.

A selection of available map types are:

- BPF_ARRAY,

- BPF_PERCPU_ARRAY
- BPF_HASH
- BPF_PERCPU_HASH
- BPF_STACK_TRACE
- BPF_ARRAY_OF_MAPS

The map types used in this research are the `BPF_PERCPU_ARRAY`, `BPF_HASH`, and `BPF_PERCPU_HASH`.

The Array type has a predefined number of elements and is available per CPU or shared between cores. The maximum key size is 32-bit and the size of the array is defined before compiling [12].

The Hash type is a HashMap (or associative array) implementation which is optimised for fast lookup and updates. The Hash type uses the `jhash` library [12], which is a Linux kernel library. The maximum key size is configurable, though the standard is 64-bit.

To summarise, there are three components that enable the agnostic augmentation of DNS software. Firstly, the extension of cBPF allows BPF programs to be used in a wider method. Secondly, the three checks of the BPF verifier ensure that a BPF program is executed safely and that it is finite. Lastly, BPF maps offer a method for BPF programs to interact with user space and need to be configured with a specific type at compile time.

4 Methods

In this section we discuss the approaches for designing the prototypes. Code for all prototypes is found in the Appendix.

As mentioned before, we choose two main use cases to write prototypes for: rewriting query names of DNS messages and response rate limiting large volumes query sources. The selection of prototypes is chosen because each one demonstrates one or more different key functionalities of XDP BPF.

The QName rewrite prototype demonstrates the ability to rewrite packet contents before they are passed on to the network stack and the usability of XDP programs in large scale network tests. The general RRL prototype demonstrates a performance comparison of a native DNS software functionality, and displays the interaction with the map user space storage from within the XDP program. The per IP RRL prototype iterates on the general RRL prototype and attempts to show that increasing the number of source IPs does not influence the workload, which displays that the XDP program scales for more IPs as used as source addresses. The unknown host prototype attempts to demonstrate a functionality that is not available in most, current DNS software and shows interaction and configurability of a running XDP program from user space.

The listed functionalities help to answer the first research sub-question: Which features from XDP BPF could be used to augment DNS software?

4.1 QName rewriting

With the QName rewrite prototype we explore rewriting incoming packets by rewriting the packet before it is passed on to the network stack. This process is completely invisible to programs running in user space. The use case here is DNS message size measurements done with RIPE Atlas⁴. For this use case only the QName and checksum are rewritten, so that the query invokes a different response. The need for rewriting the query comes from the fact that the measurement has been scheduled with predefined properties allowing it to target all resolvers on all current and future probes, in an hourly schedule, and in an ongoing manner without preset end time. However, these predefined properties are bound to a fixed query name. To invoke different response sizes on different moments in time the query name needs to be changed before it is passed on to the network stack. Furthermore, the query structure involves the probe ID, which can be up to 7 digits and changes per probe. Since different probes send queries with different sizes, they invoke different respective responses. To enable all probes to invoke the same response, which is desirable for the use case, the number of digits in the ID is compensated by the XDP program. To guarantee the same sized responses for every probe ID, the query name needs to be further adjusted to reflect the length of the probe ID which is compensated with a tuned response by the NSD daemon serving the request in user space.

To this end, a DNS request is received on the network interface and the XDP program is run by the kernel. The program determines if the packet is either IPv4 or IPv6, and verifies it as a DNS message by looking at the UDP port. When the message is verified to be DNS, the expected size of the query and the size of the probe ID are checked, and a following label is expanded to fill the discrepancy between the expected and the current size. To guarantee the checksum is still correct, it is recalculated and updated. When this succeeds the packet is passed along to the network stack.

Using this method, the DNS label can be modified on the machine running the NSD instance, without the need to modify any zones or the program itself.

4.2 Response Rate Limiting

To explore the RRL topic we present three prototype iterations:

- General RRL
- Per IP RRL
- Unknown host RRL

The first prototype version, general RRL, offers an operator the ability to set a fine-grained limit of received packets that the DNS service can process without dropping packets. This limit can be fine-tuned to the capabilities of the machine and network that this service is running on.

The general version of the RRL prototype functions by counting every DNS message that is received within a predetermined time frame, and stores the

⁴The RIPE Atlas internet measurement network: <https://atlas.ripe.net>

frame starting time and the total counted packets in the frame in an array BPF map with a single entry to be accessible at every program execution. The time frame determines the granularity of the algorithm and is configured before compiling. The starting time in the frame is updated every instance it has reached or surpassed the predetermined time frame size value. This check is performed by subtracting the current time of the packet from the starting time, both in nanoseconds, and examining this against the time frame threshold. If the total number of packets exceeds the threshold within the time frame, which are both configured before compiling, all other packets in the same time frame are either dropped or replied to the sender with the truncated flag set (bounced), depending on the desired behavior.

Per IP RRL

The per IP RRL prototype functions much the same as the general version. The main difference is that the time frame and number of packets are calculated and stored per IP. The values are stored in a hash type BPF map where the key is the source IP address of the incoming query. The size of the map can be configured before compilation.

The per IP RRL prototype provides the option to differentiate between low volume query sources and high volume, repetitive query sources. This differentiation allows the algorithm to restrict high volume senders, while serving the low volumes senders unimpeded.

Unknown host RRL

The unknown host RRL prototype is similar to the per IP RRL prototype in that IPs can be rate limited, but differs in how this is accomplished.

The unknown host RRL program functions by matching the IP of the current DNS message to an entry in the known hosts HashMap, and subjects the message to RRL if the source IP is not in the allowlist. If an entry exists in the known hosts HashMap for an IP, the program can be configured to bounce or pass the message based on the entry. The entry, like in the per IP RRL prototype, stores the number of packets seen from this sender and the start of the current time frame. If the entry does not exist, the algorithm creates an entry in a second HashMap with the number of packets set to 1 and the start time set to the current time. This second HashMap is identical in use to the one in the per IP RRL program.

The method of adding entries to the HashMap is manually adding them from user space. This method allows the user to create a list of trusted senders which are exempt from rate limiting without recompiling and reloading the BPF program.

5 Experiments

In this Section we describe the experimental setups for the two use cases, QName rewriting and RRL, and their respective prototypes.

5.1 QName rewrite

The QName rewrite prototype is used in research by Koolhaas and Slokker [2]. The goal of their experiment is to find the optimal size for User Datagram Protocol (UDP) packets traversing the network without fragmenting using the RIPE ATLAS measurement platform by measuring DNS responses of a known size.

Their experimental setup involves long running RIPE Atlas measurements for which the parameters cannot be changed once the measurements are started. The query name (QName) of DNS requests is changed by an XDP program to invoke requests for different response sizes. Rewriting by the XDP program is done without the need to change the query parameters of the RIPE Atlas measurement, such as individual probe IDs. The setup used by Koolhaas and Slokker can be seen in Figure 3

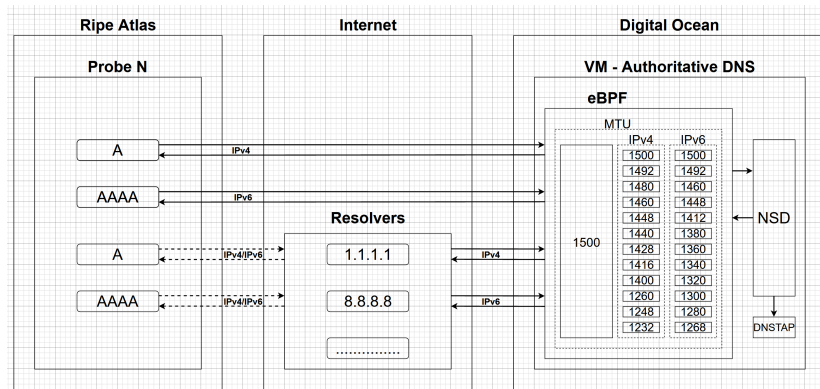


Figure 3: The experimental setup used by Koolhaas and Slokker. Note that the MTU size is determined by the query and is corrected in the eBPF program. (copied from [2]).

An example query from probe N could request a response of 1480 bytes directly from the authoritative server. In this example case, the ID of probe N, which is included in the query, is smaller than the maximum probe length and the response the query receives would be 1472 bytes. To ensure each measurement receives the same response size regardless of the ID length, the eBPF program adds padding to the query to ensure the response size is correct.

The QName prototype is used extensively during measurements of the Koolhaas and Slokker research and shown to function for all used queries which are also shown in Figure 3.

5.2 General RRL

The experimental setup for the RRL experiment consists of two servers. One running an instance of NSD⁵, as DNS server, and one running an instance of the Flamethrower⁶ tool. Flamethrower is a community-based tool that can be

⁵NSD, an open source authoritative DNS name server: <https://www.nlnetlabs.nl/projects/nsd/about/>

⁶The open source Flamethrower source code: <https://github.com/DNS-OARC/flamethrower>

used as a stress test method for DNS software by sending many queries in rapid succession. The tool has the capability of concurrent requests with configurable Queries per second (QPS). Flamethrower is used in all RRL experiments.

At the time of writing, an instance of NSD comes preconfigured with the RRL feature enabled with a default of 200 QPS per IP address. Enabling the NSD RRL implementation for this experiment offers a method of verifying the RRL prototype. Most packets will be dropped by the NSD RRL when querying in large volumes. Counting the number of dropped packets, as timeouts, allows us to measure the effectiveness of the XDP program as we configure it to return the queries that are response rate limited with the truncated flag set. Using increasingly more aggressive RRL thresholds, i.e. lowering the threshold so that more messages are Response Rate Limited, in the XDP program, we expect the number of dropped packets to decrease. If successful, we can conclude that the program is functional.

Increasing the RRL threshold while measuring the CPU load will also give us insight in the efficiency relation between the two. We expect the CPU load to drop when increasing the aggressiveness of the RRL threshold on the XDP program, by lowering the RRL threshold.

The configuration choice of using the NSD RRL functionality is made, as logging the packet dropped by the XDP program would be a relatively computationally expensive operation. An expensive operation could influence the CPU load measurements, and the current solution allows for easier measuring, as Flamethrower creates a report for every executed run.

5.3 Per IP RRL

The setup for the per IP RRL experiment is similar to the previous experiment. The difference is that the RRL functionality of NSD is switched off and that instead of receiving the same query from one source IP, here we send traffic from multiple source IP addresses. Since the functionality of the program is established by the previous experiment, the goal of this experiment is assessing how traffic originating from multiple senders affects the CPU load of our XDP implementation. The NSD RRL functionality is not needed for this goal. A visual representation of the setup can be seen in Figure 4.

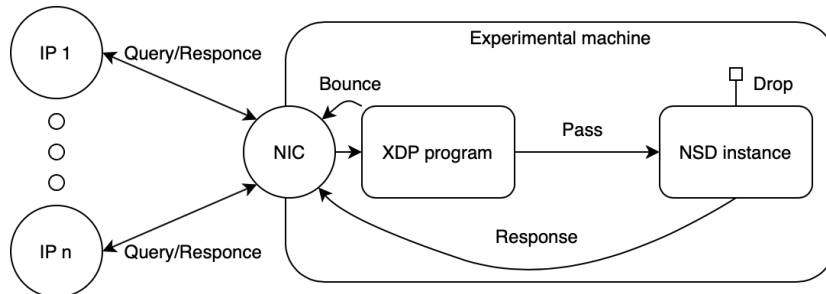


Figure 4: Representation of the setup of the per IP experiment. Note that all senders have their own assigned IP address from which they send queries to the same machine.

To keep the measurements comparable while varying the number of senders, a maximum QPS limit needs to be found. The machine running the flamethrower instances and the machine running the NSD instance are not on the same network. We empirically determine a maximum QPS for our specific measurement setup, ensuring the load per core never exceeds 100 percent usage, as this could influence the results. When found, this limit will be used in the rest of the experiments.

For this experiment we configure the RRL aggressively, to ensure the majority of the responses are sent by the XDP program, as opposed to NSD. Since the total number of responses does not vary significantly, we do not expect the CPU load of the NSD instance to change as it is highly optimised for large number of queries from different sources.

To measure the difference in CPU load from the XDP program, the number of flamethrower instances with their own IPs is increased incrementally. We expect the total CPU load not to change significantly as the number of source IPs increases.

5.4 Unknown host RRL

The experimental setup for the unknown host RRL experiment is the same as for the per IP RRL experiment. It consist of multiple flamethrower instances with their own respective IP addresses, which send queries to a single machine running an NSD instance, such as can be seen in Figure 4.

The difference compared to the previous experiment is that instead of varying the number of flamethrower instances, the number of source IPs that are included in the allowlist is incremented between runs. This incremental change allows us to verify that the XDP program is functional, as we expect to see CPU load increase as more IPs are included in the allowlist as the workload for NSD grows. IPs are added to the map of known hosts with the `bpftool` toolset.

Querying non-existing domains allows us to differentiate the responses received from the XDP program and NSD. The responses that the Flamethrower tools receive contain the DNS response code `REFUSED`, while responses that are above the RRL threshold within the XDP program contain the DNS return code `NOERROR`. This difference in response code allows us to verify that the number of queries that reach the NSD instance is increasing according to the respective run.

6 Results

In this Section we present our results from developing and assessing the prototypes as described in Section 5.

6.1 QName rewrite

With the research of Koolhaas and Slokker, a RIPE Atlas measurement was scheduled with the special properties of reaching all resolvers running on all RIPE Atlas probes, and going on continuously and persistently. These special

properties came at the expense of the ability to modify the query. Our query rewrite XDP program resolved this limitation independent of, and without modifying DNS software. Maintaining standard components greatly increases reproducibility of the experiment. Furthermore, we were able to fine tune rewritten queries on a per packet basis to fine tune the variations in response sizes caused by the length of the Probe ID embedded in the query name.

We found that rewriting the query name of incoming DNS packets, and restoring the original query name on response, is possible as DNS wire format imposes hard limits on the length and the number of the labels in a query name, which makes the number of control flow paths and boundaries within the program, that the verifier has to check, finite. Overall, we find that an entire DNS packet can be rewritten as long as the size of the packet does not change. A change of the size of a packet is limited to 256 bytes. These 256 bytes are located at the head of the packet to enable encapsulation and are less suitable for extension of the DNS packet [13].

The functionality of the QName rewrite program is emergent from the results of [2], as they verify that the program is effective. All queries within specifications are handled correctly by the XDP program, including the layer 3 and layer 4 checksums. Queries outside of the specification are not tested, and would not likely result in correct execution. Koolhaas and Slokker have reported no loss in performance using the XDP program. From their report we can conclude that it is functional in large scale networking experiments, and thus realises the goal set for this experiment in Section 4.

6.2 General RRL

As mentioned in Section 5, we perform multiple iterations of the experiment increasing the aggressiveness of the RRL threshold to verify that the general RRL program works correctly. To this end, we choose 8 RRL thresholds which give insight into the performance of the XDP program. This performance is measured by comparing timeouts, i.e. packets that are dropped by NSD, against truncated packets. The range of these thresholds is found by experimentally finding a low threshold number that results in a number of timeouts that is close to zero. From there, the threshold is increased in intermediate steps to a number that is close to the number of packets that are dropped by the NSD resolver itself. To visualise the working of the program, we configure all packets that go above the threshold to be returned with the truncation flag instead of being dropped.

In total, roughly 600.000 DNS queries are sent per run. We find that with this configuration NSD uses roughly 80% of a single CPU core for replies. Since the assigned CPU core changes per run, we measure the total CPU load. The machine running the Flamethrower instance has access to 3 CPU cores.

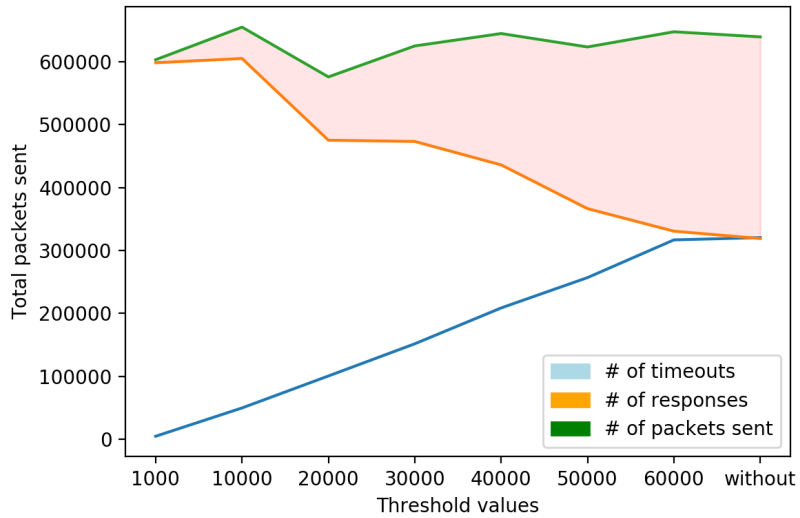


Figure 5: Visualisation of the general RRL experiment. Note that the red shaded area shows all the packets that timeout at the Flamethrower tool as they are dropped by the built-in NSD RRL, and the final measurement is taken without an XDP program attached. We discuss the variability of the number of packets sent in Section 7.

Figure 5 shows incrementally increases in the configured RRL threshold until it roughly corresponds with the NSD RRL threshold. From the figure we notice that increasing the RRL threshold, decreases the number of responses that are received by the Flamethrower tool. This is due to more queries reaching the NSD instance, which starts dropping them accordingly. From this observation we can deduce that packets are processed correctly by the XDP program and conclude that the program is functional.

Simultaneously to the measurements, the total CPU load on the resolver machine was measured as well.

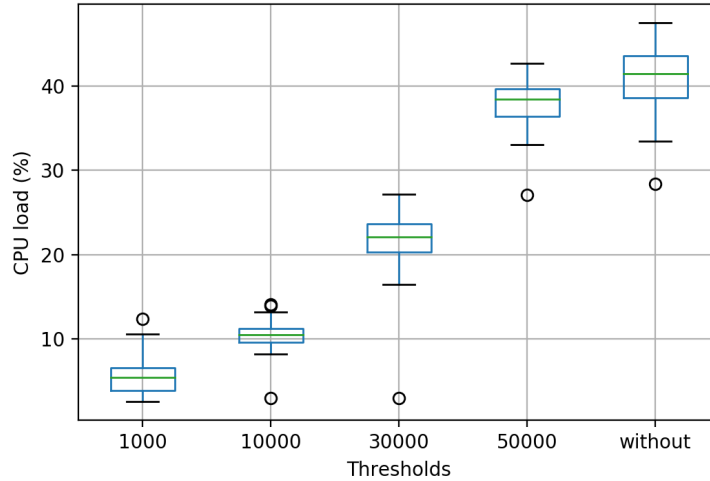


Figure 6: Visualisation of the measured total CPU load per configured RRL threshold. Note that this change is due to a difference in workload for NSD.

In Figure 6 we visualise the CPU load of different RRL thresholds. We can observe that the CPU load drops significantly when lowering the RRL threshold. When we compare the XDP program to the NSD RRL functionality, we can observe that the CPU load is lower in the XDP program while handling the same number of packets. We can observe that this behavior increases as the XDP RRL is configured more aggressively.

From this observation we can conclude that using the XDP program is more CPU efficient than using the NSD RRL functionality. This conclusion is in line with our expectations from Section 5.2.

6.3 Per IP RRL

The maximum QPS described in Section 5.3 to keep the total CPU load around 80% without dropping any packets is found to be roughly 45.000 QPS. All individual Flamethrower processes are assigned one IP address and limit their QPS per source IP so that the total QPS never rises above 45.000 QPS. For example, in the configuration with 2 source IPs, both flamethrower instances are configured to limit at 22.500 QPS. The configured RRL threshold in the XDP program used for all tests is 1 QPS, which is the minimum in the program. This low threshold ensures that we only measure the change in the XDP workload and not in a significantly changing NSD workload. The total CPU load is measured every second for 30 seconds per run.

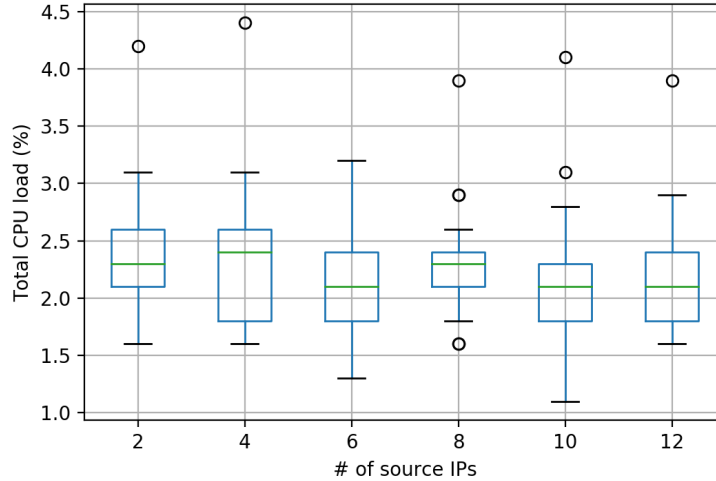


Figure 7: Visualisation of the measured total CPU load per number of source IPs. Note that the QPS is constant over the experiment, and the RRL threshold for the XDP program is set to 1 QPS per source IP.

In Figure 7 we can observe the measured total CPU load for the different number of source IP addresses. We notice the medians of all measurements to be within 0.5% of each other, so we can conclude that our predictions made in Section 5.3 are correct and the total CPU load does not increase significantly as the number of source IPs increases.

6.4 Unknown host RRL

The source IPs are added to the allowlist in increments of 2 IPs per run. This increment gives us a gradual overview of the change in CPU load as more messages reach the NSD instance. Each run involves 12 instances of Flamethrower which are configured to a limit of 3750 QPS. The RRL threshold of the XDP program is configured at 1000 QPS to make a change in CPU load per run visible.

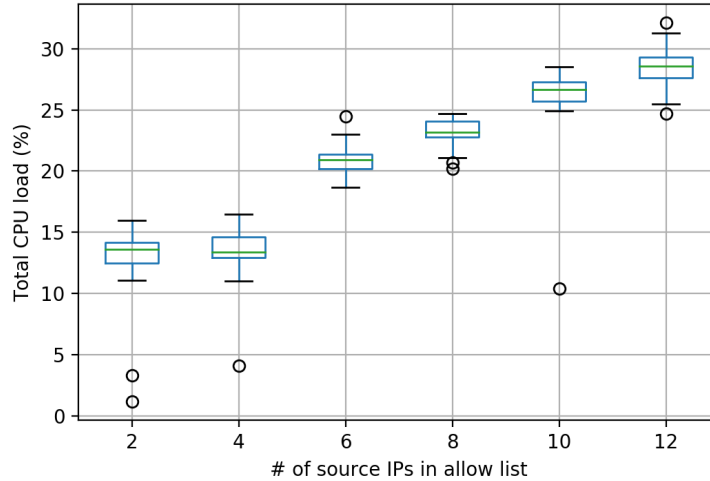


Figure 8: A visualisation of the measured total CPU load per number of source IPs in the allowlist. Note that number of IPs used as source address and the QPS is constant throughout the experiment and only the IPs added to the allowlist changes per run.

From Figure 8 we can observe the visualisation of the CPU load per number of IPs used as source address in the allowlist. The visualisation shows that adding IPs to the allowlist increases the CPU load. Because the number of source IPs and the QPS limit is constant throughout the experiment, we can conclude that the increase in CPU load is due to NSD handling the queries, and the XDP program is functional.

7 Discussion

In the results from Figure 5 we can observe that the number of packets sent is not the same for every data point. While these measurements are taken over 10 second intervals, as mentioned before in Section 5.3 the information and control of the network falls outside of the resources of this research. While the network variability does influence the results, we argue that this is not significant for the conclusions of the results, as the variability is never more than 5 percent.

While this research shows that XDP BPF provides an opportunity to augment DNS services without making changes to the service itself, the augmentation shown is limited to DNS over UDP.

Currently DNS over UDP is the standard and possible agnostic augmentations can therefore be implemented in the foreseeable future. We do note that DNS over TLS or DNS over HTTPS specifications exist and acknowledge that the class of augmentations shown in this work would not be functional, as the decryption of packets happens in higher network layers which negates the usability on lower layers.

In all RRL experiments, every packet is checked if it exceeds the configured RRL threshold in the time frame, and if it arrived within the time frame. While this method of checking is functional, it could be argued that this is not efficient as BPF helper calls, such as the helper call for the current time, are relatively computationally expensive as they make a system call. To optimise this, the time check could be configured based on a predetermined number of packets or a per-packet probability. This change can be useful in high traffic volume environments, as the expensive system call and the check could be superfluous for every packet, as large volumes of packets are received every time frame.

While these checking methods are likely an efficiency improvement over checking every packet, such performance improvements are not within the scope of exploring this research

In the general RRL experiments the XDP program is configured to bounce the received DNS queries with the TC bit set if the RRL threshold is exceeded. This behavior is not strictly RRL, as this requires the packet to be dropped instead of bounced. The XDP program contrasts the behavior of the NSD instance, which by default drops 50% of the requests when they exceed the configured NSD RRL threshold. While bouncing packets is computationally more expensive as it requires more steps than dropping the packet, the conclusion drawn in Section 6.2 is still valid, as both actions happen at the kernel level before the network stack and the results show the XDP program to be more efficient than the NSD functionality, even while bouncing the packet instead of dropping it.

A possible criticism on the proposed method of augmenting DNS software is that it results in layer violation of the application layer in the OSI model. While we concede this point as application layer data is resolved on the wrong layer, we can also argue that the violation is contained to a single machine running the program and is thus constrained to the machine and not the wider network.

8 Conclusion

In this research, we have proposed a novel method of augmenting DNS services agnostic of the DNS software supplier. This method is enabled by the BPF instruction set and the included XDP kernel hook.

To showcase the different aspects of this technology we have created 4 prototypes. With the QName rewrite prototype we have shown the ability to rewrite packet contents before they are passed on to the network stack. With the general RRL experiments we have drawn a performance comparison between a DNS software native functionality and the general RRL prototype, and has displayed interaction with the map user space storage from within the XDP program. The per IP RRL prototype iterates on the general RRL prototype and shows that increasing the number of source IPs does not influence the workload. The unknown host prototype has demonstrated a functionality that is not available in most current DNS software and has shown configurability of a running XDP program from user space.

We have answered the first research sub-question in Section 3: Which features from XDP eBPF could be used to augment DNS software? We answer this by examining the properties of BPF and XDP.

The implementations described in Section 4 answer the second research sub-question: How can DNS augmentations be implemented based upon these XDP eBPF features? We answer this with the created prototypes.

In Section 6 we have concluded that all experiments show the XDP BPF prototypes lowers the CPU load compared to the user space counterpart, and have therefore answered the third research sub-question: How do these implementations impact performance?

With this research we have shown that XDP BPF is a viable candidate for augmentations to DNS software agnostic of the supplier and answered the main research question: How can XDP BPF be used to augment and improve DNS software?

9 Future work

To further this research, experiments around the per IP RRL prototype could be done on high traffic simulations or high volume name servers. This could provide a more real-world scenario where DNS defense against amplification attacks is required, while legitimate traffic could still be resolved.

Another potential candidate for augmenting a DNS service with BPF, is server-side handling of DNS Cookies. DNS Cookies are an in-protocol security mechanism against off path denial-of-service and amplification, forgery, or cache poisoning attacks. This protection comes from the DNS server recognising a returning client by verifying a cookie carried in the request, generated from a client-provided nonce and a secret known only by the server.

The DNS Cookie interactions are embedded in EDNS0 options and piggy back on, and are completely independent from, regular DNS interactions between a client and a server, which makes them especially suitable for processing by BPF. Furthermore, earlier DNS software implementations were not interoperable making them impractical to deploy on multi-vendor anycast networks. An BPF implementation would resolve this issue by providing a single, consistent and performant implementation of DNS Cookie processing regardless of the DNS software used to serve the DNS requests.

Another interesting topic for future research is to examine the possible security implications of the prototypes presented in this research. While they are bound to the limitations and the verifier of BPF, they could still pose a possible security risk as unwanted instruction executions in kernel space could pose a security risk.

References

- [1] O. Hohlfeld, J. Krude, J. H. Reelfs, J. R uth, and K. Wehrle, “Demystifying the Performance of XDP BPF,” in *2019 IEEE Conference on Network Softwarization (NetSoft)*. IEEE, 2019, pp. 208–212.

- [2] A. Koolhaas and T. Slokker, “Defragmenting DNS; Determining the optimal maximum UDP response size for DNS,” 2020.
- [3] S. Goldlust, “A Quick Introduction to Response Rate Limiting,” <https://kb.isc.org/docs/aa-01000> (Accessed: 15 July 2020), 2018.
- [4] M. A. Vieira, M. S. Castanho, R. D. Pacífico, E. R. Santos, E. P. C. Júnior, and L. F. Vieira, “Fast Packet Processing with eBPF and XDP: Concepts, Code, Challenges, and Applications,” *ACM Computing Surveys (CSUR)*, vol. 53, no. 1, pp. 1–36, 2020.
- [5] T. Høiland-Jørgensen, J. D. Brouer, D. Borkmann, J. Fastabend, T. Herbert, D. Ahern, and D. Miller, “The eXpress Data Path: Fast programmable packet processing in the operating system kernel,” in *Proceedings of the 14th International Conference on Emerging Networking Experiments and Technologies*, 2018, pp. 54–66.
- [6] H. Wieren, “Signature-Based DDoS Attack Mitigation: Automated Generating Rules for Extended Berkeley Packet Filter and Express Data Path,” Master’s thesis, University of Twente, 2019.
- [7] N. de Bruijn, “eBPF Based Networking,” 2017.
- [8] G. Bertin, “XDP in practice: integrating XDP into our DDoS mitigation pipeline,” in *Technical Conference on Linux Networking, Netdev*, vol. 2, 2017.
- [9] A. Fabre, “L4drop: XDP DDOS mitigations,” <https://blog.cloudflare.com/l4drop-xdp-ebpf-based-ddos-mitigations/> (Accessed: Jun 2020), 2019.
- [10] B. Gregg, “Learn ebpf tracing: Tutorial and examples,” <http://www.brendangregg.com/blog/2019-01-01/learn-ebpf-tracing.html> (Accessed: Aug 2020), 1 Jan 2019.
- [11] M. Fleming, “A thorough introduction to eBPF,” *Linux Weekly News*, 2017.
- [12] C. Neira, “bcc reference guide,” https://github.com/iovisor/bcc/blob/master/docs/reference_guide.md#1-bpf_table (Accessed: Aug 2020), 2020.
- [13] Cilium, “BPF and XDP Reference Guide,” <https://docs.cilium.io/en/v1.8/bpf/#xdp> (Accessed: Aug 2020), 2020.

10 Acknowledgements

I would like to thank my supervisors Willem Toorop and Luuk Hendriks for giving me the opportunity to work on NLnet Labs Research on Networks project, the great cooperation, and their invaluable contributions to this research.

I would also like to thank my family and friends for the (continued) support before, during, and after this project.

A QName rewrite prototype code

```
1  /*
2  *  rrl-per-ip
3  *  Implements per IP RLL within a time frame for hosts that are not known in hte
4  *  Jun 2020 - Tom Carpay
5  */
6
7  #include <stdint.h>
8  #include <linux/bpf.h>
9  #include <linux/if_ether.h> /* for struct ethhdr */
10 #include <linux/ip.h>      /* for struct iphdr */
11 #include <linux/ipv6.h>    /* for struct ipv6hdr */
12 #include <linux/in.h>      /* for IPPROTO_UDP */
13 #include <linux/udp.h>     /* for struct udphdr */
14 #include <linux/pkt_cls.h>
15 #include <bpf_helpers.h>
16
17 #define DNS_PORT 53
18 #define MAX_LABELS 50
19
20 #ifndef __section
21 # define __section(NAME)          \
22     __attribute__((section(NAME), used))
23 #endif
24
25 #ifndef __inline
26 # define __inline                \
27     inline __attribute__((always_inline))
28 #endif
29
30 #if __BYTE_ORDER__ == __ORDER_LITTLE_ENDIAN__
31 # ifndef ntohs
32 #  define ntohs(x) __builtin_bswap16(x)
33 # endif
34 # ifndef htons
35 #  define htons(x) __builtin_bswap16(x)
36 # endif
37 # ifndef ntohl
38 #  define ntohl(x) __builtin_bswap32(x)
39 # endif
40 # ifndef htonl
41 #  define htonl(x) __builtin_bswap32(x)
42 # endif
43 #else
44 # ifndef ntohs
45 #  define ntohs(x) (x)
46 # endif
47 # ifndef htons
48 #  define htons(x) (x)
```

```

49  # endif
50  # ifndef ntohl
51  # define ntohl(x) (x)
52  # endif
53  # ifndef htonl
54  # define htonl(x) (x)
55  # endif
56  #endif
57
58  #ifndef memset
59  # define memset(dest, chr, n)  __builtin_memset((dest), (chr), (n))
60  #endif
61
62  #ifndef memcpy
63  # define memcpy(dest, src, n)  __builtin_memcpy((dest), (src), (n))
64  #endif
65
66  #ifndef memmove
67  # define memmove(dest, src, n) __builtin_memmove((dest), (src), (n))
68  #endif
69
70
71  struct vlanhdr {
72      uint16_t tci;
73      uint16_t encap_proto;
74  };
75
76  struct dnshdr {
77      uint16_t id;
78
79      uint8_t  rd      : 1;
80      uint8_t  tc      : 1;
81      uint8_t  aa      : 1;
82      uint8_t  opcode  : 4;
83      uint8_t  qr      : 1;
84
85      uint8_t  rcode   : 4;
86      uint8_t  cd      : 1;
87      uint8_t  ad      : 1;
88      uint8_t  z       : 1;
89      uint8_t  ra      : 1;
90
91      uint16_t qdcount;
92      uint16_t ancount;
93      uint16_t nscount;
94      uint16_t arcount;
95  };
96
97  struct cursor {
98      void *pos;

```



```

99         void *end;
100     };
101
102     static __inline
103     void cursor_init(struct cursor *c, struct xdp_md *ctx)
104     {
105         c->end = (void *) (long) ctx->data_end;
106         c->pos = (void *) (long) ctx->data;
107     }
108
109     static __inline
110     void cursor_init_skb(struct cursor *c, struct __sk_buff *skb)
111     {
112         c->end = (void *) (long) skb->data_end;
113         c->pos = (void *) (long) skb->data;
114     }
115
116     #define PARSE_FUNC_DECLARATION(STRUCT)
117     static __inline
118     struct STRUCT *parse_ ## STRUCT (struct cursor *c)
119     {
120         struct STRUCT *ret = c->pos;
121         if (c->pos + sizeof(struct STRUCT) > c->end)
122             return 0;
123         c->pos += sizeof(struct STRUCT);
124         return ret;
125     }
126
127     PARSE_FUNC_DECLARATION(ethhdr)
128     PARSE_FUNC_DECLARATION(vlanhdr)
129     PARSE_FUNC_DECLARATION(iphdr)
130     PARSE_FUNC_DECLARATION(ipv6hdr)
131     PARSE_FUNC_DECLARATION(udphdr)
132     PARSE_FUNC_DECLARATION(dnshdr)
133
134     static __inline
135     struct ethhdr *parse_eth(struct cursor *c, uint16_t *eth_proto)
136     {
137         struct ethhdr *eth;
138
139         if (!(eth = parse_ethhdr(c)))
140             return 0;
141
142         *eth_proto = eth->h_proto;
143         if (*eth_proto == htons(ETH_P_8021Q)
144             || *eth_proto == htons(ETH_P_8021AD)) {
145             struct vlanhdr *vlan;
146
147             if (!(vlan = parse_vlanhdr(c)))
148                 return 0;

```

```

149
150     *eth_proto = vlan->encap_proto;
151     if (*eth_proto == htons(ETH_P_8021Q)
152         || *eth_proto == htons(ETH_P_8021AD)) {
153         if (!(vlan = parse_vlanhdr(c)))
154             return 0;
155
156         *eth_proto = vlan->encap_proto;
157     }
158 }
159 return eth;
160 }
161
162 static __inline
163 void update_checksum(uint16_t *csum, uint16_t old_val, uint16_t new_val)
164 {
165     uint32_t new_csum_value;
166     uint32_t new_csum_comp;
167     uint32_t undo;
168
169     undo = ~((uint32_t)*csum) + ~((uint32_t)old_val);
170     new_csum_value = undo + (undo < ~((uint32_t)old_val)) + (uint32_t)new_val;
171     new_csum_comp = new_csum_value + (new_csum_value < ((uint32_t)new_val));
172     new_csum_comp = (new_csum_comp & 0xFFFF) + (new_csum_comp >> 16);
173     new_csum_comp = (new_csum_comp & 0xFFFF) + (new_csum_comp >> 16);
174     *csum = (uint16_t)~new_csum_comp;
175 }
176
177 static __inline
178 void rewrite_qname4(struct cursor *c, uint8_t *pkt, struct udphdr *udp)
179 {
180     uint8_t *labels[MAX_LABELS];
181     uint8_t i;
182
183     for (i = 0; i < MAX_LABELS; i++) { /* Maximum 128 labels */
184         uint8_t o;
185
186         if (c->pos + 1 > c->end)
187             return;
188
189         o = *(uint8_t *)c->pos;
190         if ((o & 0xC0) == 0xC0) {
191             return;
192
193         } else if (o & 0xC0)
194             /* Unknown label type */
195             return;
196
197         labels[i] = c->pos;
198         c->pos += o + 1;

```

```

199         if (!o)
200             break;
201     }
202     if (i >= MAX_LABELS || i < 5
203         || *labels[i-4] != 10
204         || labels[i-4] + *labels[i-4] + 2 > (uint8_t *)c->end
205         || labels[i-4][ 1] < '0' || labels[i-4][1] > '9'
206         || labels[i-4][ 2] < '0' || labels[i-4][2] > '9'
207         || labels[i-4][ 3] < '0' || labels[i-4][3] > '9'
208         || labels[i-4][ 4] < '0' || labels[i-4][4] > '9'
209         || labels[i-4][ 5] != '-')
210         || (labels[i-4][ 6] != 'p' && labels[i-4][6] != 'P')
211         || (labels[i-4][ 7] != 'l' && labels[i-4][7] != 'L')
212         || (labels[i-4][ 8] != 'u' && labels[i-4][8] != 'U')
213         || (labels[i-4][ 9] != 's' && labels[i-4][9] != 'S')
214         || labels[i-4][10] != '0' )
215         return;
216
217     /* Change aligned on 16 bits for checksum recalculation */
218     uint16_t *pls_pos = (labels[i-4] + 10 - (uint8_t *)udp) % 2
219         ? (uint16_t *)&labels[i-4][9]
220         : (uint16_t *)&labels[i-4][10];
221     uint16_t old_pls = *pls_pos;
222
223     switch (*labels[i-5]) {
224     case 39: labels[i-4][10] = '2'; break;
225     case 38: labels[i-4][10] = '4'; break;
226     case 37: labels[i-4][10] = '6'; break;
227     case 36: labels[i-4][10] = '8'; break;
228     case 35: labels[i-4][10] = 'a'; break;
229     case 34: labels[i-4][10] = 'c'; break;
230     default: break;
231     }
232     update_checksum(&udp->check, old_pls, *pls_pos);
233     #if MTU4 != 1500
234     if (labels[i-4][1] != '1' || labels[i-4][2] != '5'
235         || labels[i-4][3] != '0' || labels[i-4][4] != '0')
236         return;
237
238     if ((labels[i-4] - (uint8_t *)udp) % 2) {
239         uint16_t *sh1_pos = (uint16_t *)&labels[i-4][1];
240         uint16_t old_sh1 = *sh1_pos;
241         uint16_t *sh2_pos = (uint16_t *)&labels[i-4][3];
242         uint16_t old_sh2 = *sh2_pos;
243
244         labels[i-4][1] = MTU4_STR[0];
245         labels[i-4][2] = MTU4_STR[1];
246         labels[i-4][3] = MTU4_STR[2];
247         labels[i-4][4] = MTU4_STR[3];
248

```

```

249         update_checksum(&udp->check, old_sh1,*sh1_pos);
250         update_checksum(&udp->check, old_sh2,*sh2_pos);
251     } else {
252         uint16_t *sh1_pos = (uint16_t*)&labels[i-4][0];
253         uint16_t old_sh1 = *sh1_pos;
254         uint16_t *sh2_pos = (uint16_t*)&labels[i-4][2];
255         uint16_t old_sh2 = *sh2_pos;
256         uint16_t *sh3_pos = (uint16_t*)&labels[i-4][4];
257         uint16_t old_sh3 = *sh3_pos;
258
259         labels[i-4][1] = MTU4_STR[0];
260         labels[i-4][2] = MTU4_STR[1];
261         labels[i-4][3] = MTU4_STR[2];
262         labels[i-4][4] = MTU4_STR[3];
263
264         update_checksum(&udp->check, old_sh1,*sh1_pos);
265         update_checksum(&udp->check, old_sh2,*sh2_pos);
266         update_checksum(&udp->check, old_sh3,*sh3_pos);
267     }
268     #endif
269 }
270
271 static __inline
272 void restore_qname4(struct cursor *c, uint8_t *pkt, struct udphdr *udp)
273 {
274     uint8_t *labels[MAX_LABELS];
275     uint8_t i;
276
277     for (i = 0; i < MAX_LABELS; i++) { /* Maximum 128 labels */
278         uint8_t o;
279
280         if (c->pos + 1 > c->end)
281             return;
282
283         o = *(uint8_t *)c->pos;
284         if ((o & 0xC0) == 0xC0) {
285             return;
286
287         } else if (o & 0xC0)
288             /* Unknown label type */
289             return;
290
291         labels[i] = c->pos;
292         c->pos += o + 1;
293         if (!o)
294             break;
295     }
296     if (i >= MAX_LABELS || i < 5
297         || *labels[i-4] != 10
298         || labels[i-4] + *labels[i-4] + 2 > (uint8_t *)c->end

```

```

299     || labels[i-4][ 1] < '0' || labels[i-4][1] > '9'
300     || labels[i-4][ 2] < '0' || labels[i-4][2] > '9'
301     || labels[i-4][ 3] < '0' || labels[i-4][3] > '9'
302     || labels[i-4][ 4] < '0' || labels[i-4][4] > '9'
303     || labels[i-4][ 5] != '-'
304     || (labels[i-4][ 6] != 'p' && labels[i-4][6] != 'P')
305     || (labels[i-4][ 7] != 'l' && labels[i-4][7] != 'L')
306     || (labels[i-4][ 8] != 'u' && labels[i-4][8] != 'U')
307     || (labels[i-4][ 9] != 's' && labels[i-4][9] != 'S'))
308         return;
309
310     if (labels[i-4][10] != '0') {
311         /* Change aligned on 16 bits for checksum recalculation
312          * Doesn't work on TC/TX! Maybe we should use bpf_l4_csum_replace()
313          * and bpf_csum_diff().
314          *
315          * uint16_t *pls_pos = (labels[i-4] + 10 - (uint8_t *)udp) % 2
316          *                    ? (uint16_t *)&labels[i-4][9]
317          *                    : (uint16_t *)&labels[i-4][10];
318          * uint16_t old_pls = *pls_pos;
319          */
320         labels[i-4][10] = '0';
321         // udp->check = 0;
322     }
323
324     #if MTU4 != 1500
325     if (labels[i-4][1] != MTU4_STR[0] || labels[i-4][2] != MTU4_STR[1]
326         || labels[i-4][3] != MTU4_STR[2] || labels[i-4][4] != MTU4_STR[3])
327         return;
328
329     if ((labels[i-4] - (uint8_t *)udp) % 2) {
330         /* TODO: 4 bytes checksum recalculating labels[i-4][1-4]
331          *        with bpf_l4_csum_replace() and bpf_csum_diff()
332          */
333         labels[i-4][1] = '1';
334         labels[i-4][2] = '5';
335         labels[i-4][3] = '0';
336         labels[i-4][4] = '0';
337         // udp->check = 0;
338     } else {
339         /* TODO: 6 bytes checksum recalculating labels[i-4][0-5]
340          *        with bpf_l4_csum_replace() and bpf_csum_diff()
341          */
342         labels[i-4][1] = '1';
343         labels[i-4][2] = '5';
344         labels[i-4][3] = '0';
345         labels[i-4][4] = '0';
346         // udp->check = 0;
347     }
348     #endif

```

```

349         return;
350     }
351
352     static __inline
353     void rewrite_qname6(struct cursor *c, uint8_t *pkt, struct udphdr *udp)
354     {
355         uint8_t *labels[MAX_LABELS];
356         uint8_t i;
357
358         for (i = 0; i < MAX_LABELS; i++) { /* Maximum 128 labels */
359             uint8_t o;
360
361             if (c->pos + 1 > c->end)
362                 return;
363
364             o = *(uint8_t *)c->pos;
365             if ((o & 0xC0) == 0xC0) {
366                 return;
367
368             } else if (o & 0xC0)
369                 /* Unknown label type */
370                 return;
371
372             labels[i] = c->pos;
373             c->pos += o + 1;
374             if (!o)
375                 break;
376         }
377         if (i >= MAX_LABELS || i < 5
378             || *labels[i-4] != 10
379             || labels[i-4] + *labels[i-4] + 2 > (uint8_t *)c->end
380             || labels[i-4][ 1] < '0' || labels[i-4][1] > '9'
381             || labels[i-4][ 2] < '0' || labels[i-4][2] > '9'
382             || labels[i-4][ 3] < '0' || labels[i-4][3] > '9'
383             || labels[i-4][ 4] < '0' || labels[i-4][4] > '9'
384             || labels[i-4][ 5] != '-'
385             || (labels[i-4][ 6] != 'p' && labels[i-4][6] != 'P')
386             || (labels[i-4][ 7] != 'l' && labels[i-4][7] != 'L')
387             || (labels[i-4][ 8] != 'u' && labels[i-4][8] != 'U')
388             || (labels[i-4][ 9] != 's' && labels[i-4][9] != 'S')
389             || labels[i-4][10] != '0' )
390             return;
391
392         /* Change aligned on 16 bits for checksum recalculaction */
393         uint16_t *pls_pos = (labels[i-4] + 10 - (uint8_t *)udp) % 2
394             ? (uint16_t *)&labels[i-4][9]
395             : (uint16_t *)&labels[i-4][10];
396         uint16_t old_pls = *pls_pos;
397
398         switch (*labels[i-5]) {

```

```

399     case 39: labels[i-4][10] = '2'; break;
400     case 38: labels[i-4][10] = '4'; break;
401     case 37: labels[i-4][10] = '6'; break;
402     case 36: labels[i-4][10] = '8'; break;
403     case 35: labels[i-4][10] = 'a'; break;
404     case 34: labels[i-4][10] = 'c'; break;
405     default: break;
406     }
407     update_checksum(&udp->check, old_pls, *pls_pos);
408     #if MTU6 != 1500
409     if (labels[i-4][1] != '1' || labels[i-4][2] != '5'
410         || labels[i-4][3] != '0' || labels[i-4][4] != '0')
411         return;
412
413     if ((labels[i-4] - (uint8_t *)udp) % 2) {
414         uint16_t *sh1_pos = (uint16_t*)&labels[i-4][1];
415         uint16_t old_sh1 = *sh1_pos;
416         uint16_t *sh2_pos = (uint16_t*)&labels[i-4][3];
417         uint16_t old_sh2 = *sh2_pos;
418
419         labels[i-4][1] = MTU6_STR[0];
420         labels[i-4][2] = MTU6_STR[1];
421         labels[i-4][3] = MTU6_STR[2];
422         labels[i-4][4] = MTU6_STR[3];
423
424         update_checksum(&udp->check, old_sh1, *sh1_pos);
425         update_checksum(&udp->check, old_sh2, *sh2_pos);
426     } else {
427         uint16_t *sh1_pos = (uint16_t*)&labels[i-4][0];
428         uint16_t old_sh1 = *sh1_pos;
429         uint16_t *sh2_pos = (uint16_t*)&labels[i-4][2];
430         uint16_t old_sh2 = *sh2_pos;
431         uint16_t *sh3_pos = (uint16_t*)&labels[i-4][4];
432         uint16_t old_sh3 = *sh3_pos;
433
434         labels[i-4][1] = MTU6_STR[0];
435         labels[i-4][2] = MTU6_STR[1];
436         labels[i-4][3] = MTU6_STR[2];
437         labels[i-4][4] = MTU6_STR[3];
438
439         update_checksum(&udp->check, old_sh1, *sh1_pos);
440         update_checksum(&udp->check, old_sh2, *sh2_pos);
441         update_checksum(&udp->check, old_sh3, *sh3_pos);
442     }
443     #endif
444 }
445
446 static __inline
447 uint16_t restore_qname6(struct cursor *c, uint8_t *pkt, struct udphdr *udp)
448 {

```

```

449     uint8_t *labels[MAX_LABELS];
450     uint8_t i;
451
452     for (i = 0; i < MAX_LABELS; i++) { /* Maximum 128 labels */
453         uint8_t o;
454
455         if (c->pos + 1 > c->end)
456             return 0;
457
458         o = *(uint8_t *)c->pos;
459         if ((o & 0xC0) == 0xC0) {
460             return 0;
461
462         } else if (o & 0xC0)
463             /* Unknown label type */
464             return 0;
465
466         labels[i] = c->pos;
467         c->pos += o + 1;
468         if (!o)
469             break;
470     }
471     if (i >= MAX_LABELS || i < 5
472         || *labels[i-4] != 10
473         || labels[i-4] + *labels[i-4] + 2 > (uint8_t *)c->end
474         || labels[i-4][ 1] < '0' || labels[i-4][1] > '9'
475         || labels[i-4][ 2] < '0' || labels[i-4][2] > '9'
476         || labels[i-4][ 3] < '0' || labels[i-4][3] > '9'
477         || labels[i-4][ 4] < '0' || labels[i-4][4] > '9'
478         || labels[i-4][ 5] != '-'
479         || (labels[i-4][ 6] != 'p' && labels[i-4][6] != 'P')
480         || (labels[i-4][ 7] != 'l' && labels[i-4][7] != 'L')
481         || (labels[i-4][ 8] != 'u' && labels[i-4][8] != 'U')
482         || (labels[i-4][ 9] != 's' && labels[i-4][9] != 'S'))
483         return 0;
484
485     if (labels[i-4][10] != '0') {
486         if ((labels[i-4] - (uint8_t *)udp) % 2) {
487             uint16_t old = ((uint16_t *)&labels[i-4][1])[4];
488             labels[i-4][10] = '0';
489             #if MTU6 > 1280
490                 update_checksum(&udp->check, old, ((uint16_t *)&labels[i-4][1])[4]);
491             #endif
492         } else {
493             uint16_t old = ((uint16_t *)labels[i-4])[5];
494             labels[i-4][10] = '0';
495             #if MTU6 > 1280
496                 update_checksum(&udp->check, old, ((uint16_t *)labels[i-4])[5]);
497             #endif
498         }

```



```

499     }
500
501     #if MTU6 != 1500
502         if (labels[i-4][1] != MTU6_STR[0] || labels[i-4][2] != MTU6_STR[1]
503             || labels[i-4][3] != MTU6_STR[2] || labels[i-4][4] != MTU6_STR[3])
504             return 0;
505
506         if ((labels[i-4] - (uint8_t *)udp) % 2) {
507             uint16_t old0 = ((uint16_t *)&labels[i-4][1])[0];
508             uint16_t old1 = ((uint16_t *)&labels[i-4][1])[1];
509             labels[i-4][1] = '1';
510             labels[i-4][2] = '5';
511             labels[i-4][3] = '0';
512             labels[i-4][4] = '0';
513         #if MTU6 > 1280
514             update_checksum(&udp->check, old0, ((uint16_t *)&labels[i-4][1])[0]);
515             update_checksum(&udp->check, old1, ((uint16_t *)&labels[i-4][1])[1]);
516         #endif
517         } else {
518             uint16_t old0 = ((uint16_t *)labels[i-4])[0];
519             uint16_t old1 = ((uint16_t *)labels[i-4])[1];
520             uint16_t old2 = ((uint16_t *)labels[i-4])[2];
521             labels[i-4][1] = '1';
522             labels[i-4][2] = '5';
523             labels[i-4][3] = '0';
524             labels[i-4][4] = '0';
525         #if MTU6 > 1280
526             update_checksum(&udp->check, old0, ((uint16_t *)labels[i-4])[0]);
527             update_checksum(&udp->check, old1, ((uint16_t *)labels[i-4])[1]);
528             update_checksum(&udp->check, old2, ((uint16_t *)labels[i-4])[2]);
529         #endif
530         }
531     #endif
532     return (labels[i-4] - (uint8_t *)pkt) + 1;
533 }
534
535
536 __section("xdp-rewrite-qname")
537 int xdp_rewrite_qname(struct xdp_md *ctx)
538 {
539     struct cursor    c;
540     uint16_t         eth_proto;
541     struct iphdr     *ipv4;
542     struct ipv6hdr   *ipv6;
543     struct udphdr    *udp;
544     struct dnshdr    *dns;
545
546     cursor_init(&c, ctx);
547     if (!parse_eth(&c, &eth_proto))
548         return XDP_PASS;

```

```

549
550     if (eth_proto == htons(ETH_P_IP)) {
551         if (!(ipv4 = parse_iphdr(&c)) || ipv4->protocol != IPPROTO_UDP
552             || !(udp = parse_udphdr(&c)) || udp->dest != htons(DNS_PORT)
553             || !(dns = parse_dnshdr(&c)))
554             return XDP_PASS;
555
556         rewrite_qname4(&c, (void *)dns, udp);
557
558     } else if (eth_proto == htons(ETH_P_IPV6)) {
559         if (!(ipv6 = parse_ipv6hdr(&c)) || ipv6->nexthdr != IPPROTO_UDP
560             || !(udp = parse_udphdr(&c)) || udp->dest != htons(DNS_PORT)
561             || !(dns = parse_dnshdr(&c)))
562             return XDP_PASS;
563
564         rewrite_qname6(&c, (void *)dns, udp);
565     }
566     return XDP_PASS;
567 }
568
569 __section("tc-restore-qname")
570 int tc_restore_qname(struct __sk_buff *skb)
571 {
572     struct cursor    c;
573     uint16_t         eth_proto;
574     struct iphdr     *ipv4;
575     struct ipv6hdr   *ipv6;
576     struct udphdr    *udp;
577     struct dnshdr    *dns;
578
579     cursor_init_skb(&c, skb);
580     if (!parse_eth(&c, &eth_proto))
581         return TC_ACT_OK;
582
583     if (eth_proto == htons(ETH_P_IP)) {
584         if (!(ipv4 = parse_iphdr(&c)) || ipv4->protocol != IPPROTO_UDP
585             || !(udp = parse_udphdr(&c)) || udp->source != htons(DNS_PORT)
586             || !(dns = parse_dnshdr(&c)))
587             return TC_ACT_OK;
588
589         uint16_t old_val = ipv4->frag_off;
590         #if __BYTE_ORDER__ == __ORDER_LITTLE_ENDIAN__
591             ipv4->frag_off |= 0x0040;
592         #else
593             ipv4->frag_off |= 0x4000;
594         #endif
595         update_checksum(&ipv4->check, old_val, ipv4->frag_off);
596         restore_qname4(&c, (void *)dns, udp);
597     } else if (eth_proto == htons(ETH_P_IPV6)) {

```

```
599         if (!(ipv6 = parse_ipv6hdr(&c)) || ipv6->nexthdr != IPPROTO_UDP
600             || !(udp = parse_udphdr(&c)) || udp->source != htons(DNS_PORT)
601             || !(dns = parse_dnshdr(&c)))
602             return TC_ACT_OK;
603
604         restore_qname6(&c, (void *)dns, udp);
605     }
606     return TC_ACT_OK;
607 }
608
609 char __license[] __section("license") = "GPL";
```

B General RLL prototype code

```
1  /*
2  *  General_RRL.c
3  *  Implements a semi fine grained udp_dns_reply RRL within a time frame
4  *  Jun 2020 - Tom Carpay
5  */
6
7  /*
8  *  Includes
9  */
10 #include <stdint.h>
11 #include <linux/bpf.h>
12 #include <bpf_helpers.h>    /* for bpf_get_prandom_u32() */
13 #include <bpf_endian.h>    /* for __bpf_htons() */
14 #include <linux/if_ether.h> /* for struct ethhdr */
15 #include <linux/ip.h>      /* for struct iphdr */
16 #include <linux/ipv6.h>    /* for struct ipv6hdr */
17 #include <linux/in.h>      /* for IPPROTO_UDP */
18 #include <linux/udp.h>     /* for struct udphdr */
19 #include <string.h>        /* for memcpy() */
20
21 /*
22 *  Begin defines
23 */
24 #define DNS_PORT    53
25
26 #define FRAME_SIZE    1000000000
27 #define THRESHOLD    50000
28 /*
29 *  End defines
30 */
31
32 /*
33 **  Store the time frame
34 */
35 struct bucket {
36     uint64_t start_time;
37     uint64_t n_packets;
38     // uint64_t qps;
39 };
40
41 struct bpf_map_def SEC("maps") state_map = {
42     .type = BPF_MAP_TYPE_PERCPU_ARRAY,
43     .key_size = sizeof(uint32_t),
44     .value_size = sizeof(struct bucket),
45     .max_entries = 1
46 };
47
48 /*
```

```

49  * Store the VLAN header
50  */
51  struct vlanhdr {
52      uint16_t tci;
53      uint16_t encap_proto;
54  };
55
56  /*
57  * Store the DNS header
58  */
59  struct dnshdr {
60      uint16_t id;
61      union {
62          struct {
63              uint8_t rd      : 1;
64              uint8_t tc      : 1;
65              uint8_t aa      : 1;
66              uint8_t opcode  : 4;
67              uint8_t qr      : 1;
68
69              uint8_t rcode   : 4;
70              uint8_t cd      : 1;
71              uint8_t ad      : 1;
72              uint8_t z        : 1;
73              uint8_t ra      : 1;
74          } as_bits_and_pieces;
75          uint16_t as_value;
76      } flags;
77      uint16_t qdcount;
78      uint16_t ancount;
79      uint16_t nscount;
80      uint16_t arcount;
81  };
82
83  /*
84  * Helper pointer to parse the incoming packets
85  */
86  struct cursor {
87      void *pos;
88      void *end;
89  };
90
91
92  /*
93  * Initializer of a cursor pointer
94  */
95  static __always_inline
96  void cursor_init(struct cursor *c, struct xdp_md *ctx)
97  {
98      c->end = (void *) (long) ctx->data_end;

```

```

99         c->pos = (void *) (long) ctx->data;
100     }
101
102     #define PARSE_FUNC_DECLARATION(STRUCT)
103     static __always_inline
104     struct STRUCT *parse_ ## STRUCT (struct cursor *c)
105     {
106         struct STRUCT *ret = c->pos;
107         if (c->pos + sizeof(struct STRUCT) > c->end)
108             return 0;
109         c->pos += sizeof(struct STRUCT);
110         return ret;
111     }
112
113     PARSE_FUNC_DECLARATION(ethhdr)
114     PARSE_FUNC_DECLARATION(vlanhdr)
115     PARSE_FUNC_DECLARATION(iphdr)
116     PARSE_FUNC_DECLARATION(ipv6hdr)
117     PARSE_FUNC_DECLARATION(udphdr)
118     PARSE_FUNC_DECLARATION(dnshdr)
119
120     /*
121      * Parse ethernet frame and fill the struct
122      */
123     static __always_inline
124     struct ethhdr *parse_eth(struct cursor *c, uint16_t *eth_proto)
125     {
126         struct ethhdr *eth;
127
128         if (!(eth = parse_ethhdr(c)))
129             return 0;
130
131         *eth_proto = eth->h_proto;
132         if (*eth_proto == __bpf_htons(ETH_P_8021Q)
133             || *eth_proto == __bpf_htons(ETH_P_8021AD)) {
134             struct vlanhdr *vlan;
135
136             if (!(vlan = parse_vlanhdr(c)))
137                 return 0;
138
139             *eth_proto = vlan->encap_proto;
140             if (*eth_proto == __bpf_htons(ETH_P_8021Q)
141                 || *eth_proto == __bpf_htons(ETH_P_8021AD)) {
142                 if (!(vlan = parse_vlanhdr(c)))
143                     return 0;
144
145                 *eth_proto = vlan->encap_proto;
146             }
147         }
148         return eth;

```

```

149 }
150
151 /*
152  * Recalculate the checksum
153  */
154 static __always_inline
155 void update_checksum(uint16_t *csum, uint16_t old_val, uint16_t new_val)
156 {
157     uint32_t new_csum_value;
158     uint32_t new_csum_comp;
159     uint32_t undo;
160
161     undo = ~((uint32_t)*csum) + ~((uint32_t)old_val);
162     new_csum_value = undo + (undo < ~((uint32_t)old_val)) + (uint32_t)new_val;
163     new_csum_comp = new_csum_value + (new_csum_value < ((uint32_t)new_val));
164     new_csum_comp = (new_csum_comp & 0xFFFF) + (new_csum_comp >> 16);
165     new_csum_comp = (new_csum_comp & 0xFFFF) + (new_csum_comp >> 16);
166     *csum = (uint16_t)~new_csum_comp;
167 }
168
169 /*
170  * Parse DNS message.
171  * Returns 1 if message needs to go through (i.e. pass)
172  *        -1 if something went wrong and the packet needs to be dropped
173  *        0 if (modified) message needs to be replied
174  */
175 static __always_inline
176 int udp_dns_reply(struct cursor *c)
177 {
178     struct udphdr *udp;
179     struct dnshdr *dns;
180     uint32_t key;
181
182     // check that we have a DNS packet
183     if (!(udp = parse_udphdr(c)) || udp->dest != __bpf_htons(DNS_PORT)
184         || !(dns = parse_dnshdr(c)))
185         return 1;
186
187     // get the starting time frame from the map
188     key = 0;
189     struct bucket *b = bpf_map_lookup_elem(&state_map, &key);
190
191     // the bucket must exist
192     if (!b)
193     {
194         //bpf_printk("!FRAME \n");
195         return -1;
196     }
197
198     // increment number of packets

```

```

199     b->n_packets++;
200
201
202     // @TODO evaluate this option for frame timing
203     // look at the timing every 100 packets
204     if (b->n_packets % 100 == 0)
205
206     // look at the timing of the packet a percentage of the time
207     //if (bpf_get_prandom_u32() % 100 < 50)
208     {
209         // get the current and elapsed time
210         uint64_t now = bpf_ktime_get_ns();
211         uint64_t elapsed = now - b->start_time;
212
213         // make sure the elapsed time is set and not outside of the frame
214         if (b->start_time == 0 || elapsed >= FRAME_SIZE)
215         {
216             //bpf_printk("New timeframe\n");
217             // start new time frame
218             b->start_time = now;
219             b->n_packets = 0;
220         }
221     }
222
223     //bpf_printk("n_packets: %llu\n", b->n_packets);
224
225     // @TODO refine bounce rate to fit curve
226     if (b->n_packets > THRESHOLD)
227     {
228         //bpf_printk("bounce\n");
229         //save the old header values
230         uint16_t old_val = dns->flags.as_value;
231
232         // change the DNS flags
233         dns->flags.as_bits_and_pieces.ad = 0;
234         dns->flags.as_bits_and_pieces.qr = 1;
235         dns->flags.as_bits_and_pieces.tc = 1;
236
237         // change the UDP destination to the source
238         udp->dest = udp->source;
239         udp->source = __bpf_htons(DNS_PORT);
240
241         // calculate and write the new checksum
242         update_checksum(&udp->check, old_val, dns->flags.as_value);
243
244         // bounce
245         return 0;
246     }
247     else
248     {

```



```

249             // pass
250             return 1;
251         }
252     }
253
254     /*
255     * Recieve and parse request
256     * @var struct xdp_md
257     */
258     SEC("xdp-dns-too-many")
259     int xdp_dns_too_many(struct xdp_md *ctx)
260     {
261         // store variables
262         struct cursor    c;
263         struct ethhdr    *eth;
264         uint16_t         eth_proto;
265         struct iphdr     *ipv4;
266         struct ipv6hdr   *ipv6;
267         int              r = 0;
268
269         // initialise the cursor
270         cursor_init(&c, ctx);
271         if (!(eth = parse_eth(&c, &eth_proto)))
272             return XDP_PASS;
273
274         // differentiate the parsing of the IP header based on the version
275         if (eth_proto == __bpf_htons(ETH_P_IP)) {
276             if (!(ipv4 = parse_iphdr(&c))
277                 || ipv4->protocol != IPPROTO_UDP
278                 || (r = udp_dns_reply(&c))) {
279                 return r < 0 ? XDP_ABORTED : XDP_PASS;
280             }
281
282             uint32_t swap_ipv4 = ipv4->daddr;
283             ipv4->daddr = ipv4->saddr;
284             ipv4->saddr = swap_ipv4;
285
286         } else if (eth_proto == __bpf_htons(ETH_P_IPV6)) {
287             if (!(ipv6 = parse_ipv6hdr(&c))
288                 || ipv6->nexthdr != IPPROTO_UDP
289                 || (r = udp_dns_reply(&c)))
290                 return r < 0 ? XDP_ABORTED : XDP_PASS;
291
292             struct in6_addr swap_ipv6 = ipv6->daddr;
293             ipv6->daddr = ipv6->saddr;
294             ipv6->saddr = swap_ipv6;
295         } else {
296             return XDP_PASS;
297         }
298

```

```
299     uint8_t swap_eth[ETH_ALEN];
300     memcpy(swap_eth, eth->h_dest, ETH_ALEN);
301     memcpy(eth->h_dest, eth->h_source, ETH_ALEN);
302     memcpy(eth->h_source, swap_eth, ETH_ALEN);
303
304     // bounce the request
305     return XDP_TX;
306 }
307
308 char __license[] SEC("license") = "GPL";
```

C Per IP RRL prototype code

```
1  /*
2  *  rrl-per-ip
3  *  Implements a semi fine grained udp_dns_reply RRL per ip address within a time frame
4  *  Jun 2020 - Tom Carpay
5  */
6
7  /*
8  *  Includes
9  */
10 #include <stdint.h>
11 #include <linux/bpf.h>
12 #include <bpf_helpers.h>    /* for bpf_get_prandom_u32() */
13 #include <bpf_endian.h>    /* for __bpf_htons() */
14 #include <linux/if_ether.h> /* for struct ethhdr */
15 #include <linux/ip.h>      /* for struct iphdr */
16 #include <linux/ipv6.h>    /* for struct ipv6hdr */
17 #include <linux/in.h>      /* for IPPROTO_UDP */
18 #include <linux/udp.h>     /* for struct udphdr */
19 #include <string.h>        /* for memcpy() */
20
21 /*
22 *  Begin defines
23 */
24 #define DNS_PORT      53
25
26 #define FRAME_SIZE    1000000000
27 #define THRESHOLD     1000
28 /*
29 *  End defines
30 */
31
32 /*
33 *  Store the time frame
34 */
35 struct bucket {
36     uint64_t start_time;
37     uint64_t n_packets;
38 };
39
40 struct bpf_map_def SEC("maps") state_map = {
41     .type = BPF_MAP_TYPE_PERCPU_HASH,
42     .key_size = sizeof(uint32_t),
43     .value_size = sizeof(struct bucket),
44     .max_entries = 100
45 };
46
47 struct bpf_map_def SEC("maps") state_map_v6 = {
48     .type = BPF_MAP_TYPE_PERCPU_HASH,
```

```

49     .key_size = sizeof(struct in6_addr),
50     .value_size = sizeof(struct bucket),
51     .max_entries = 100
52 };
53
54 /*
55  * Store the VLAN header
56  */
57 struct vlanhdr {
58     uint16_t tci;
59     uint16_t encap_proto;
60 };
61
62 /*
63  * Store the DNS header
64  */
65 struct dnshdr {
66     uint16_t id;
67     union {
68         struct {
69             uint8_t rd    : 1;
70             uint8_t tc    : 1;
71             uint8_t aa    : 1;
72             uint8_t opcode : 4;
73             uint8_t qr    : 1;
74
75             uint8_t rcode  : 4;
76             uint8_t cd    : 1;
77             uint8_t ad    : 1;
78             uint8_t z     : 1;
79             uint8_t ra    : 1;
80         } as_bits_and_pieces;
81         uint16_t as_value;
82     } flags;
83     uint16_t qdcount;
84     uint16_t ancourt;
85     uint16_t nscount;
86     uint16_t arcount;
87 };
88
89 /*
90  * Helper pointer to parse the incoming packets
91  */
92 struct cursor {
93     void *pos;
94     void *end;
95 };
96
97
98 /*

```

```

99  * Initializer of a cursor pointer
100 */
101 static __always_inline
102 void cursor_init(struct cursor *c, struct xdp_md *ctx)
103 {
104     c->end = (void *) (long) ctx->data_end;
105     c->pos = (void *) (long) ctx->data;
106 }
107
108 #define PARSE_FUNC_DECLARATION(STRUCT)
109 static __always_inline
110 struct STRUCT *parse_## STRUCT (struct cursor *c)
111 {
112     struct STRUCT *ret = c->pos;
113     if (c->pos + sizeof(struct STRUCT) > c->end)
114         return 0;
115     c->pos += sizeof(struct STRUCT);
116     return ret;
117 }
118
119 PARSE_FUNC_DECLARATION(ethhdr)
120 PARSE_FUNC_DECLARATION(vlanhdr)
121 PARSE_FUNC_DECLARATION(iphdr)
122 PARSE_FUNC_DECLARATION(ipv6hdr)
123 PARSE_FUNC_DECLARATION(udphdr)
124 PARSE_FUNC_DECLARATION(dnshdr)
125
126 /*
127  * Parse ethernet frame and fill the struct
128  */
129 static __always_inline
130 struct ethhdr *parse_eth(struct cursor *c, uint16_t *eth_proto)
131 {
132     struct ethhdr *eth;
133
134     if (!(eth = parse_ethhdr(c)))
135         return 0;
136
137     *eth_proto = eth->h_proto;
138     if (*eth_proto == __bpf_htons(ETH_P_8021Q)
139         || *eth_proto == __bpf_htons(ETH_P_8021AD)) {
140         struct vlanhdr *vlan;
141
142         if (!(vlan = parse_vlanhdr(c)))
143             return 0;
144
145         *eth_proto = vlan->encap_proto;
146         if (*eth_proto == __bpf_htons(ETH_P_8021Q)
147             || *eth_proto == __bpf_htons(ETH_P_8021AD)) {
148             if (!(vlan = parse_vlanhdr(c)))

```

```

149             return 0;
150
151             *eth_proto = vlan->encap_proto;
152         }
153     }
154     return eth;
155 }
156
157 /*
158  * Recalculate the checksum
159  */
160 static __always_inline
161 void update_checksum(uint16_t *csum, uint16_t old_val, uint16_t new_val)
162 {
163     uint32_t new_csum_value;
164     uint32_t new_csum_comp;
165     uint32_t undo;
166
167     undo = ~((uint32_t)*csum) + ~((uint32_t)old_val);
168     new_csum_value = undo + (undo < ~((uint32_t)old_val)) + (uint32_t)new_val;
169     new_csum_comp = new_csum_value + (new_csum_value < ((uint32_t)new_val));
170     new_csum_comp = (new_csum_comp & 0xFFFF) + (new_csum_comp >> 16);
171     new_csum_comp = (new_csum_comp & 0xFFFF) + (new_csum_comp >> 16);
172     *csum = (uint16_t)~new_csum_comp;
173 }
174
175 static __always_inline
176 int do_rate_limit(struct udphdr *udp, struct dnshdr *dns, struct bucket *b)
177 {
178     // increment number of packets
179     b->n_packets++;
180
181     // get the current and elapsed time
182     uint64_t now = bpf_ktime_get_ns();
183     uint64_t elapsed = now - b->start_time;
184
185     // make sure the elapsed time is set and not outside of the frame
186     if (b->start_time == 0 || elapsed >= FRAME_SIZE)
187     {
188         //bpf_printk("New timeframe\n");
189         // start new time frame
190         b->start_time = now;
191         b->n_packets = 0;
192     }
193
194     // @TODO refine bounce rate to fit curve
195     if (b->n_packets < THRESHOLD)
196         return 1;
197
198     //bpf_printk("bounce\n");

```

```

199     //save the old header values
200     uint16_t old_val = dns->flags.as_value;
201
202     // change the DNS flags
203     dns->flags.as_bits_and_pieces.ad = 0;
204     dns->flags.as_bits_and_pieces.qr = 1;
205     dns->flags.as_bits_and_pieces.tc = 1;
206
207     // change the UDP destination to the source
208     udp->dest = udp->source;
209     udp->source = __bpf_htons(DNS_PORT);
210
211     // calculate and write the new checksum
212     update_checksum(&udp->check, old_val, dns->flags.as_value);
213
214     // bounce
215     return 0;
216 }
217
218 /*
219  * Parse DNS ipv4 message
220  * Returns 1 if message needs to go through (i.e. pass)
221  *        -1 if something went wrong and the packet needs to be dropped
222  *        0 if (modified) message needs to be replied
223  */
224 static __always_inline
225 int udp_dns_reply_v4(struct cursor *c, uint32_t key)
226 {
227     struct udphdr *udp;
228     struct dns_hdr *dns;
229
230     // check that we have a DNS packet
231     if (!(udp = parse_udphdr(c)) || udp->dest != __bpf_htons(DNS_PORT)
232         || !(dns = parse_dns_hdr(c)))
233         return 1;
234
235     // get the starting time frame from the map
236     struct bucket *b = bpf_map_lookup_elem(&state_map, &key);
237
238     // the bucket must exist
239     if (b)
240         return do_rate_limit(udp, dns, b);
241
242     // create new starting bucket for this key
243     struct bucket new_bucket;
244     new_bucket.start_time = bpf_ktime_get_ns();
245     new_bucket.n_packets = 0;
246
247     // store the bucket and pass the packet
248     bpf_map_update_elem(&state_map, &key, &new_bucket, BPF_ANY);

```

```

249         return 1;
250     }
251
252     /*
253     * Parse DNS message.
254     * Returns 1 if message needs to go through (i.e. pass)
255     *     -1 if something went wrong and the packet needs to be dropped
256     *     0 if (modified) message needs to be replied
257     */
258     static __always_inline
259     int udp_dns_reply_v6(struct cursor *c, struct in6_addr *key)
260     {
261         struct udphdr *udp;
262         struct dnshdr *dns;
263
264         // check that we have a DNS packet
265         if (!(udp = parse_udphdr(c)) || udp->dest != __bpf_htons(DNS_PORT)
266             || !(dns = parse_dnshdr(c)))
267             return 1;
268
269         // get the starting time frame from the map
270         struct bucket *b = bpf_map_lookup_elem(&state_map_v6, key);
271
272         // the bucket must exist
273         if (b)
274             return do_rate_limit(udp, dns, b);
275
276         // create new starting bucket for this key
277         struct bucket new_bucket;
278         new_bucket.start_time = bpf_ktime_get_ns();
279         new_bucket.n_packets = 0;
280
281         // store the bucket and pass the packet
282         bpf_map_update_elem(&state_map_v6, key, &new_bucket, BPF_ANY);
283         return 1;
284     }
285
286
287     /*
288     * Recieve and parse request
289     * @var struct xdp_md
290     */
291     SEC("xdp-dns-too-many")
292     int xdp_dns_too_many(struct xdp_md *ctx)
293     {
294         // store variables
295         struct cursor c;
296         struct ethhdr *eth;
297         uint16_t eth_proto;
298         struct iphdr *ipv4;

```



```

299     struct ipv6hdr *ipv6;
300     int             r = 0;
301
302     // initialise the cursor
303     cursor_init(&c, ctx);
304     if (!(eth = parse_eth(&c, &eth_proto)))
305         return XDP_PASS;
306
307     // differentiate the parsing of the IP header based on the version
308     if (eth_proto == __bpf_htons(ETH_P_IP))
309     {
310         if (!(ipv4 = parse_iphdr(&c))
311             || ipv4->protocol != IPPROTO_UDP
312             || (r = udp_dns_reply_v4(&c, ipv4->saddr))) {
313
314             return r < 0 ? XDP_ABORTED : XDP_PASS;
315         }
316
317         uint32_t swap_ipv4 = ipv4->daddr;
318         ipv4->daddr = ipv4->saddr;
319         ipv4->saddr = swap_ipv4;
320
321     }
322     else if (eth_proto == __bpf_htons(ETH_P_IPV6))
323     {
324         if (!(ipv6 = parse_ipv6hdr(&c))
325             || ipv6->nexthdr != IPPROTO_UDP
326             || (r = udp_dns_reply_v6(&c, &ipv6->saddr)))
327             return r < 0 ? XDP_ABORTED : XDP_PASS;
328
329         struct in6_addr swap_ipv6 = ipv6->daddr;
330         ipv6->daddr = ipv6->saddr;
331         ipv6->saddr = swap_ipv6;
332     }
333     else
334     {
335         return XDP_PASS;
336     }
337
338     uint8_t swap_eth[ETH_ALEN];
339     memcpy(swap_eth, eth->h_dest, ETH_ALEN);
340     memcpy(eth->h_dest, eth->h_source, ETH_ALEN);
341     memcpy(eth->h_source, swap_eth, ETH_ALEN);
342
343     // Bounce
344
345     // bounce the request
346     return XDP_TX;
347 }
348

```

```
349 char __license[] SEC("license") = "GPL";
```

D Unknown host RRL prototype code

```
1  /*
2  *  rrl-per-ip
3  *  Implements per IP RLL within a time frame for hosts that are not known in the
4  *  Jun 2020 - Tom Carpay
5  */
6
7  /*
8  *  Includes
9  */
10 #include <stdint.h>
11 #include <linux/bpf.h>
12 #include <bpf_helpers.h>    /* for bpf_get_prandom_u32() */
13 #include <bpf_endian.h>    /* for __bpf_htons() */
14 #include <linux/if_ether.h> /* for struct ethhdr */
15 #include <linux/ip.h>      /* for struct iphdr */
16 #include <linux/ipv6.h>    /* for struct ipv6hdr */
17 #include <linux/in.h>      /* for IPPROTO_UDP */
18 #include <linux/udp.h>     /* for struct udphdr */
19 #include <string.h>        /* for memcpy() */
20
21 /*
22 *  Begin defines
23 */
24 #define DNS_PORT      53
25
26 #define FRAME_SIZE    1000000000
27 #define THRESHOLD     1000
28 /*
29 *  End defines
30 */
31
32 /*
33 *  Store the time frame
34 */
35 struct bucket {
36     uint64_t start_time;
37     uint64_t n_packets;
38 };
39
40 struct bpf_map_def SEC("maps") state_map = {
41     .type = BPF_MAP_TYPE_PERCPU_HASH,
42     .key_size = sizeof(uint32_t),
43     .value_size = sizeof(struct bucket),
44     .max_entries = 100 // Enough for testing purposes
45 };
46
47 struct bpf_map_def SEC("maps") state_map_v6 = {
48     .type = BPF_MAP_TYPE_PERCPU_HASH,
```

```

49     .key_size = sizeof(struct in6_addr),
50     .value_size = sizeof(struct bucket),
51     .max_entries = 100 // Enough for testing purposes
52 };
53
54 /*
55  * Smallest storage space possible
56  */
57 struct data {
58     uint8_t unused;
59 };
60
61 struct bpf_map_def SEC("maps") known_hosts = {
62     .type = BPF_MAP_TYPE_HASH,
63     .key_size = sizeof(uint32_t),
64     .value_size = sizeof(struct data),
65     .max_entries = 100 // Enough for testing purposes
66 };
67
68 struct bpf_map_def SEC("maps") known_hosts_v6 = {
69     .type = BPF_MAP_TYPE_HASH,
70     .key_size = sizeof(struct in6_addr),
71     .value_size = sizeof(struct data),
72     .max_entries = 100 // Enough for testing purposes
73 };
74
75 /*
76  * Store the VLAN header
77  */
78 struct vlanhdr {
79     uint16_t tci;
80     uint16_t encap_proto;
81 };
82
83 /*
84  * Store the DNS header
85  */
86 struct dnshdr {
87     uint16_t id;
88     union {
89         struct {
90             uint8_t rd      : 1;
91             uint8_t tc      : 1;
92             uint8_t aa      : 1;
93             uint8_t opcode  : 4;
94             uint8_t qr      : 1;
95
96             uint8_t rcode   : 4;
97             uint8_t cd      : 1;
98             uint8_t ad      : 1;

```

```

99             uint8_t z      : 1;
100             uint8_t ra     : 1;
101         }      as_bits_and_pieces;
102         uint16_t as_value;
103     } flags;
104     uint16_t qdcount;
105     uint16_t ancourt;
106     uint16_t nscourt;
107     uint16_t arcount;
108 };
109
110 /*
111  * Helper pointer to parse the incoming packets
112  */
113 struct cursor {
114     void *pos;
115     void *end;
116 };
117
118 /*
119  * Initializer of a cursor pointer
120  */
121
122 static __always_inline
123 void cursor_init(struct cursor *c, struct xdp_md *ctx)
124 {
125     c->end = (void *) (long) ctx->data_end;
126     c->pos = (void *) (long) ctx->data;
127 }
128
129 #define PARSE_FUNC_DECLARATION(STRUCT)
130 static __always_inline
131 struct STRUCT *parse_ ## STRUCT (struct cursor *c)
132 {
133     struct STRUCT *ret = c->pos;
134     if (c->pos + sizeof(struct STRUCT) > c->end)
135         return 0;
136     c->pos += sizeof(struct STRUCT);
137     return ret;
138 }
139
140 PARSE_FUNC_DECLARATION(ethhdr)
141 PARSE_FUNC_DECLARATION(vlanhdr)
142 PARSE_FUNC_DECLARATION(iphdr)
143 PARSE_FUNC_DECLARATION(ipv6hdr)
144 PARSE_FUNC_DECLARATION(udphdr)
145 PARSE_FUNC_DECLARATION(dnshdr)
146
147 /*
148  * Parse ethernet frame and fill the struct

```

```

149  */
150  static __always_inline
151  struct ethhdr *parse_eth(struct cursor *c, uint16_t *eth_proto)
152  {
153      struct ethhdr *eth;
154
155      if (!(eth = parse_ethhdr(c)))
156          return 0;
157
158      *eth_proto = eth->h_proto;
159      if (*eth_proto == __bpf_htons(ETH_P_8021Q)
160          || *eth_proto == __bpf_htons(ETH_P_8021AD)) {
161          struct vlanhdr *vlan;
162
163          if (!(vlan = parse_vlanhdr(c)))
164              return 0;
165
166          *eth_proto = vlan->encap_proto;
167          if (*eth_proto == __bpf_htons(ETH_P_8021Q)
168              || *eth_proto == __bpf_htons(ETH_P_8021AD)) {
169              if (!(vlan = parse_vlanhdr(c)))
170                  return 0;
171
172              *eth_proto = vlan->encap_proto;
173          }
174      }
175      return eth;
176  }
177
178  /*
179   * Recalculate the checksum
180   */
181  static __always_inline
182  void update_checksum(uint16_t *csum, uint16_t old_val, uint16_t new_val)
183  {
184      uint32_t new_csum_value;
185      uint32_t new_csum_comp;
186      uint32_t undo;
187
188      undo = ~((uint32_t)*csum) + ~((uint32_t)old_val);
189      new_csum_value = undo + (undo < ~((uint32_t)old_val)) + (uint32_t)new_val;
190      new_csum_comp = new_csum_value + (new_csum_value < ((uint32_t)new_val));
191      new_csum_comp = (new_csum_comp & 0xFFFF) + (new_csum_comp >> 16);
192      new_csum_comp = (new_csum_comp & 0xFFFF) + (new_csum_comp >> 16);
193      *csum = (uint16_t)~new_csum_comp;
194  }
195
196  static __always_inline
197  int do_rate_limit(struct udphdr *udp, struct dnshdr *dns, struct bucket *b)
198  {

```

```

199         // increment number of packets
200         b->n_packets++;
201
202         // get the current and elapsed time
203         uint64_t now = bpf_ktime_get_ns();
204         uint64_t elapsed = now - b->start_time;
205
206         // make sure the elapsed time is set and not outside of the frame
207         if (b->start_time == 0 || elapsed >= FRAME_SIZE)
208         {
209             //bpf_printk("New timeframe\n");
210             // start new time frame
211             b->start_time = now;
212             b->n_packets = 0;
213         }
214
215         // @TODO refine bounce rate to fit curve
216         if (b->n_packets < THRESHOLD)
217             return 1;
218
219         //bpf_printk("bounce\n");
220         //save the old header values
221         uint16_t old_val = dns->flags.as_value;
222
223         // change the DNS flags
224         dns->flags.as_bits_and_pieces.ad = 0;
225         dns->flags.as_bits_and_pieces.qr = 1;
226         dns->flags.as_bits_and_pieces.tc = 1;
227
228         // change the UDP destination to the source
229         udp->dest = udp->source;
230         udp->source = __bpf_htons(DNS_PORT);
231
232         // calculate and write the new checksum
233         update_checksum(&udp->check, old_val, dns->flags.as_value);
234
235         // bounce
236         return 0;
237     }
238
239     /*
240     * Parse DNS ipv4 message
241     * Returns 1 if message needs to go through (i.e. pass)
242     *      -1 if something went wrong and the packet needs to be dropped
243     *      0 if (modified) message needs to be replied
244     */
245     static __always_inline
246     int udp_dns_reply_v4(struct cursor *c, uint32_t key)
247     {
248         struct udphdr *udp;

```

```

249     struct dnshdr *dns;
250
251     // check that we have a DNS packet
252     if (!(udp = parse_udphdr(c)) || udp->dest != __bpf_htons(DNS_PORT)
253         || !(dns = parse_dnshdr(c)))
254         return 1;
255
256     // get the host from the list
257     struct data *host = bpf_map_lookup_elem(&known_hosts, &key);
258
259     // if the host is known we do not rate limit it
260     if (host)
261     {
262         // pass
263         return 1;
264     }
265
266     // get the starting time frame from the map
267     struct bucket *b = bpf_map_lookup_elem(&state_map, &key);
268
269
270     //bpf_printk("ip: %u\n", key);
271
272     // the bucket must exist
273     if (b)
274         return do_rate_limit(udp, dns, b);
275
276     // create new starting bucket for this key
277     struct bucket new_bucket;
278     new_bucket.start_time = bpf_ktime_get_ns();
279     new_bucket.n_packets = 0;
280
281     // store the bucket and pass the packet
282     bpf_map_update_elem(&state_map, &key, &new_bucket, BPF_ANY);
283     return 1;
284 }
285
286 /*
287  * Parse DNS message.
288  * Returns 1 if message needs to go through (i.e. pass)
289  *      -1 if something went wrong and the packet needs to be dropped
290  *      0 if (modified) message needs to be replied
291  */
292 static __always_inline
293 int udp_dns_reply_v6(struct cursor *c, struct in6_addr *key)
294 {
295     struct udphdr *udp;
296     struct dnshdr *dns;
297
298     // check that we have a DNS packet

```



```

299     if (!(udp = parse_udphdr(c)) || udp->dest != __bpf_htons(DNS_PORT)
300         || !(dns = parse_dnshdr(c)))
301         return 1;
302
303     // get the host from the list
304     struct data *host = bpf_map_lookup_elem(&known_hosts_v6, key);
305
306     // if the host is known we do not rate limit it
307     if (host)
308     {
309         // pass
310         return 1;
311     }
312
313     // get the starting time frame from the map
314     struct bucket *b = bpf_map_lookup_elem(&state_map_v6, key);
315
316     // the bucket must exist
317     if (b)
318         return do_rate_limit(udp, dns, b);
319
320     // create new starting bucket for this key
321     struct bucket new_bucket;
322     new_bucket.start_time = bpf_ktime_get_ns();
323     new_bucket.n_packets = 0;
324
325     // store the bucket and pass the packet
326     bpf_map_update_elem(&state_map_v6, key, &new_bucket, BPF_ANY);
327     return 1;
328 }
329
330
331 /*
332  * Recieve and parse request
333  * @var struct xdp_md
334  */
335 SEC("xdp-dns-too-many")
336 int xdp_dns_too_many(struct xdp_md *ctx)
337 {
338     // store variables
339     struct cursor    c;
340     struct ethhdr    *eth;
341     uint16_t         eth_proto;
342     struct iphdr     *ipv4;
343     struct ipv6hdr   *ipv6;
344     int               r = 0;
345
346     // initialise the cursor
347     cursor_init(&c, ctx);
348     if (!(eth = parse_eth(&c, &eth_proto)))

```

```

349         return XDP_PASS;
350
351         // differentiate the parsing of the IP header based on the version
352         if (eth_proto == __bpf_htons(ETH_P_IP))
353         {
354             if (!(ipv4 = parse_iphdr(&c))
355                 || ipv4->protocol != IPPROTO_UDP
356                 || (r = udp_dns_reply_v4(&c, ipv4->saddr))) {
357
358                 return r < 0 ? XDP_ABORTED : XDP_PASS;
359             }
360
361             uint32_t swap_ipv4 = ipv4->daddr;
362             ipv4->daddr = ipv4->saddr;
363             ipv4->saddr = swap_ipv4;
364
365         }
366         else if (eth_proto == __bpf_htons(ETH_P_IPV6))
367         {
368             if (!(ipv6 = parse_ipv6hdr(&c))
369                 || ipv6->nexthdr != IPPROTO_UDP
370                 || (r = udp_dns_reply_v6(&c, &ipv6->saddr)))
371                 return r < 0 ? XDP_ABORTED : XDP_PASS;
372
373             struct in6_addr swap_ipv6 = ipv6->daddr;
374             ipv6->daddr = ipv6->saddr;
375             ipv6->saddr = swap_ipv6;
376
377         }
378         else
379         {
380             return XDP_PASS;
381         }
382
383         uint8_t swap_eth[ETH_ALEN];
384         memcpy(swap_eth, eth->h_dest, ETH_ALEN);
385         memcpy(eth->h_dest, eth->h_source, ETH_ALEN);
386         memcpy(eth->h_source, swap_eth, ETH_ALEN);
387
388         // Bounce
389
390         // bounce the request
391         return XDP_TX;
392     }
393     char __license[] SEC("license") = "GPL";

```