

Leveraging eBPF for programmable network functions with IPv6 Segment Routing

Mathieu Xhonneux, Fabien Duchene, Olivier Bonaventure
 UCLouvain, Louvain-la-Neuve, Belgium
 first.last@uclouvain.be

ABSTRACT

With the advent of Software Defined Networks (SDN), Network Function Virtualisation (NFV) or Service Function Chaining (SFC), operators expect networks to support flexible services beyond the mere forwarding of packets. The network programmability framework which is being developed within the IETF by leveraging IPv6 Segment Routing enables the realisation of in-network functions.

In this paper, we demonstrate that this vision of in-network programmability can be realised. By leveraging the eBPF support in the Linux kernel, we implement a flexible framework that allows network operators to encode their own network functions as eBPF code that is automatically executed while processing specific packets. Our lab measurements indicate that the overhead of calling such eBPF functions remains acceptable. Thanks to eBPF, operators can implement a variety of network functions. We describe the architecture of our implementation in the Linux kernel. This extension has been released with Linux 4.18. We illustrate the flexibility of our approach with three different use cases: delay measurements, hybrid networks and network discovery. Our lab measurements also indicate that the performance penalty of running eBPF network functions on Linux routers does not incur a significant overhead.

CCS CONCEPTS

• **Networks** → **Programming interfaces; Middle boxes / network appliances; Programmable networks; In-network processing;**

ACM Reference Format:

Mathieu Xhonneux, Fabien Duchene, Olivier Bonaventure. 2018. Leveraging eBPF for programmable network functions with IPv6 Segment Routing. In *The 14th International Conference on emerging Networking Experiments and Technologies (CoNEXT '18), December 4–7, 2018, Heraklion, Greece*. ACM, New York, NY, USA, 6 pages. <https://doi.org/10.1145/3281411.3281426>

1 INTRODUCTION

In the late 1990s, various researchers proposed active network architectures where each packet could carry code that would be executed by intermediate routers while packets are forwarded to their final destination [42]. Several of these architectures were prototyped [40, 41] but none was deployed [16]. Still, these efforts were precursors for middleboxes [21] and Software Defined Networks (SDN) or Network Function Virtualisation (NFV). SDN enables network operators to better control the flow of packets in their networks [31, 36]. It has been mainly targeted at enterprise networks. NFV aims at enabling network operators to deploy specialised network functions which can process all the packets for specific flows. NFV is interesting for both enterprise and ISP networks. A frequently cited use case for NFV are the future 5G networks that will need to combine a variety of functions.

Segment Routing [24], initially proposed as a simplification of MultiProtocol Label Switching (MPLS) for ISP networks has gradually evolved into a much more generic solution. Segment Routing is a modern version of source routing. It enables routers to forward packets along a succession of shortest paths, each of them being identified by a segment. Besides the MPLS variant, the IETF is currently developing an IPv6 variant of Segment Routing (SRv6) that is gaining a lot of interest [17, 23]. Coupled with the fast deployment of IPv6 during the last years, this opens new opportunities for both ISP and enterprise networks.

SRv6 leverages the flexibility of the IPv6 packet format and the large IPv6 addressing space. With SRv6, each segment is encoded as an IPv6 address that is advertised through the intradomain routing protocol. NFV can be supported by assigning a specific address to each network function and using segments to forward specific flows towards the appropriate functions.

This paper focuses on the NFV use case with SRv6. We extend the SRv6 implementation in the Linux kernel [32] to support the ability to run specific network functions on a per-packet basis. To allow more flexibility compared to other solutions like P4 [15], our implementation leverages the eBPF support of the Linux kernel to enable network operators to write their own intra-domain network functions that are dynamically linked to the Linux kernel. Our evaluation shows that this enables the network functions to run efficiently inside the kernel. We then demonstrate three use cases showing very different network functions as an illustration of the flexibility of our approach.

2 BACKGROUND

Segment Routing [24] started as a modern variant of the source routing paradigm [25] using the MPLS dataplane. This architecture has now evolved to also encompass the IPv6 dataplane [38]. Segment Routing in the IPv6 data plane (SRv6) is implemented by adding an IPv6 extension header called the Segment Routing Header (SRH) [38]. This SRH contains one or more 128-bits IPv6 addresses that encode the *segments*, the nodes that must be visited on the path between the source and the destination. In the first versions of the IPv6 variant of Segment Routing, these addresses were used to identify routers and outgoing links. Several implementations of SRv6 have been announced, on commercial routers [23] and on Linux [7, 32].

In 2017, the idea of SRv6 network programming emerged [22]. It generalises the notion of segments. Each path specified by a source is decomposed into an ordered list of instructions, called *segments*. Each segment, or endpoint, represents a **function** to be executed at a specific location in the network. These functions may range from simple topological instructions (e.g. forwarding a packet on a specific link) to more complex user-defined behaviours. The SRH

carries the ordered list of segments in each packet and optionally Type-Length-Value (TLV) fields. The TLVs are 3-tuples that can be used to store additional data in the SRH, e.g. for OAM purposes.

The basic processing of packets with a SRH is the endpoint `End` function, which advances the SRH to the next segment and forwards the packet to the destination corresponding to the segment [22]. Several other functions extend this processing, such as `End.X`, which after advancing to the next segment, forwards it to a specific IPv6 nexthop, `End.T`, which performs the lookup for the next segment in a IPv6 routing table bound to the segment, `End.B6`, which inserts a new SRH on top of the existing one, etc. SRv6 actions can also be applied on packets without a SRH, e.g. inserting an SRH in an IPv6 packet, and encapsulating an outer IPv6 header with a SRH. These actions are called transit behaviours. Between segments, packets are forwarded along the shortest path.

The SRv6 implementation in the Linux kernel [7, 32] supports the basic features of SRv6 on hosts (sending and receiving IPv6 packets with an SRH) and on routers (forwarding SRv6 packets) but not all the recent SRv6 network programming features. The current Linux SRv6 implementation uses two lightweight tunnels, `seg6` and `seg6local` to support the basic SRv6 mechanisms which can be plugged in the IPv6 layer. `seg6` allows to implement the two transit behaviours mentioned above, i.e. inserting or encapsulating SRHs in traffic matching a given destination, whereas `seg6local` allows an operator to install SRv6 segments mapped to specific SRv6 functions, along with the required parameters. The set of actions provided by `seg6local` is bounded to a few simple functions statically implemented in the kernel. They do not enable extensive network programming capabilities, as required by SDN [33], SFC [6, 20] and other [9, 12, 18] applications in the SRv6 data plane.

2.1 eBPF, an in-kernel virtual machine

eBPF (for *extended Berkeley Packet Filter*), is a general-purpose virtual machine that is included in the Linux kernel since the 3.15 release. This virtual machine supports a 64 bits RISC-like CPU [3] which is an extension of the BPF virtual machine [35]. It provides a programmable interface to adapt kernel components at run-time to user-specific behaviours. While solutions such as [37, 39] use P4 [15] to achieve data plane programmability, they are limited by the fact that P4 relies on specific hardware (and/or compiler), while eBPF targets a general purpose CPU and can be used on devices like Customer-Premises Equipment (CPE). The LLVM project [5] includes a BPF backend, capable of compiling C programs to BPF bytecode. eBPF bytecode is either executed in the kernel by an interpreter or translated to native machine code using a Just-in-Time (JIT) compiler. Since the eBPF architecture is very close to the modern 64-bit ISAs, the JIT compilers usually produce efficient native code [4].

eBPF programs can be attached to predetermined hooks in the kernel. Several hooks are available in different components of the network stack, such as the traffic classifier (tc) [14], or the eXpress Data Path (XDP) [1], a low-level hook executed before the network layer, used e.g. for DDoS mitigation. When loading an eBPF program into the kernel, a verifier first ensures that it cannot threaten the stability and security of the kernel (no invalid memory accesses, possible infinite loops, ...). The eBPF program is then executed for

each packet going through the datapath associated to its hook. The program can read and, for some hooks, modify the packet.

eBPF programs can call *helper functions* [2], which are functions implemented in the kernel. They act as proxies between the kernel and the eBPF program. Using such helpers, eBPF programs can retrieve and push data from or to the kernel, and rely on mechanisms implemented in the kernel. A given hook is usually associated with a set of helpers.

There are two practical issues when developing eBPF programs. The first is how to store persistent state and the second is how it can communicate with user space applications. State can be kept persistent between multiple eBPF program invocations and shared with user space applications using *maps*. Maps are data structures implemented in the kernel as key / value stores [1]. Helpers are provided to allow eBPF programs to retrieve and store data into maps. Several structures are provided, such as arrays, hashmaps, longest prefix match tries, ... When processing packets, if information needs to be pushed asynchronously to user space, *perf events* can be used. Perf events originate from Linux's performance profiler `perf`. In a networking context, they can be used to pass custom structures from the eBPF program to the perf event ring buffer along with the packet being processed [28]. The events collected in the ring buffer can then be retrieved in user space. These mechanisms allow stateful processing and a user-space communication that would be difficult to achieve with P4.

Finally, a lightweight tunnel infrastructure named *BPF LWT*, provides generic hooks in several network layers, including IPv6 [26]. This LWT enables the execution of eBPF programs at the ingress and the egress of the routing process of network layers, but is unable to leverage the specificities of SRv6.

3 THE SRV6 EBPF INTERFACE

Many of the emerging concepts of network functions leveraging SRv6 [22] cannot be deployed using the static actions that are supported by the existing `seg6local` infrastructure on Linux [32]. Adding explicit support for each SRv6 network function in the Linux kernel would be difficult since the set of functions continues to evolve. A better approach would be to include in the Linux kernel a set of generic functions that allow network operators to implement their own SRv6 functions. For this, we propose a new eBPF interface to efficiently implement a broad range of SRv6 actions.

Our hook is a new action `End.BPF` in `seg6local`. Each instance of this action is bound to an eBPF program. It behaves as an endpoint, i.e. it only accepts SRv6 packets with a current segment corresponding to a local eBPF action, advances the SRH to the next segment, and subsequently executes the associated eBPF code. We have designed this action with two key principles in mind: (i) eBPF code cannot compromise the stability of the kernel and (ii) eBPF code should be able to leverage all the functionalities of the SRv6 data plane.

To guarantee the stability, we need to ensure that `End.BPF` can only allow write access to fields of the packet which can be modified by SRv6 endpoints. The `seg6local` actions are executed in the IPv6 layer, and further processing of the packet after `End.BPF` requires the packet to be valid. Instead of providing direct-write access to the packet to the eBPF code, we provide a specific helper function

that restrains the fields which can be modified. This function also checks whether any modification to these fields could jeopardise the integrity of the SRH.

3.1 SRv6 API for eBPF programs

All SRv6 eBPF programs are called with the packet as argument. They have full read access to its payload, starting from the outermost IPv6 header. We designed three SRv6 specific helpers to extend the functionalities of our interface:

- `bpf_lwt_seg6_store_bytes`: provides indirect write access to the editable fields of the SRH (i.e. the flags, the tag, and the TLVs).
- `bpf_lwt_seg6_adjust_srh`: allows growing or shrinking the space reserved to TLVs.
- `bpf_lwt_seg6_action`: executes a basic SRv6 function. It provides access to `End.X`, `End.T`, `End.B6`, `End.B6.Encaps` and `End.DT6`.

All SRv6 eBPF programs return an integer. This return value decides the subsequent processing of the packet. Three values can be returned:

- `BPF_OK`: a regular FIB lookup must be performed on the next segment, the packet must be forwarded on the egress interface returned by the lookup.
- `BPF_DROP`: the packet must be dropped.
- `BPF_REDIRECT`: the default endpoint lookup must not be performed, and the packet must be forwarded to the destination already set in the packet metadata.

The structure containing the packet that is passed to an eBPF program contains both the payload and metadata such as its destination. When `bpf_lwt_seg6_action` is called with an action requiring a FIB lookup (e.g. `End.X`), the helper performs the requested FIB lookup and stores the result in the metadata. When the execution of the program finishes, it is important that `End.BPF` does not execute the default lookup to the next segment afterwards, otherwise the destination previously set would be overwritten, hence the need for `BPF_REDIRECT`. If the SRH has been altered by the BPF program, a quick verification is performed to ensure that it is still valid (e.g. if the SRH has grown, ensure that the allocated space has been filled with valid TLVs), otherwise it is dropped. Finally, the packet is yielded back to the IPv6 layer, which takes care of the forwarding to the destination set in the metadata of the packet.

A fourth helper function, `bpf_lwt_push_encap`, has been implemented in the `BPF_LWT` hook. It allows to insert an SRH or encapsulate an outer IPv6 header with an SRH in pure IPv6 traffic.

Both `End.BPF` and the helpers have officially been included into the Linux kernel, and effectively released since Linux 4.18.

3.2 Performance evaluation

To measure the performance of our implementation of the `End.BPF` function, we ran several measurement campaigns in a small lab. Our lab is composed of 3 servers with Intel Xeon X3440 processors, 16GB of RAM and 10 Gbps NICs (setup 1 in Figure 1). Although these servers have multiple cores, we configured the interrupts of their NICs to direct all received packets to the same CPU core. We use `trafgen` to generate UDP packets on S1. Each UDP packet has

a payload of 64 bytes and an SRH with two segments, one bound to a function on R, and the address of S2. R executes the endpoint functions while S2 acts as a sink.

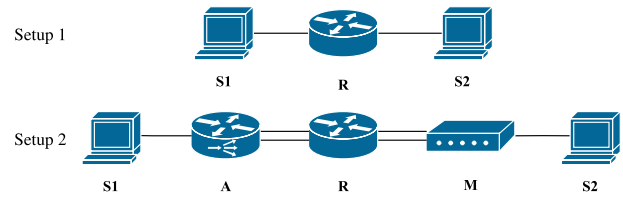


Figure 1: Lab setups used in our experiments.

For all experiments, unless stated otherwise, we enabled the JIT compiler when running our eBPF code. We first measured the raw IPv6 packet forwarding performance with those UDP packets. When the source sent 3 million packets per second, the rightmost server only received them at a rate of 610 kpps. We use this number as the reference to evaluate the impact of executing eBPF code while forwarding each packet.

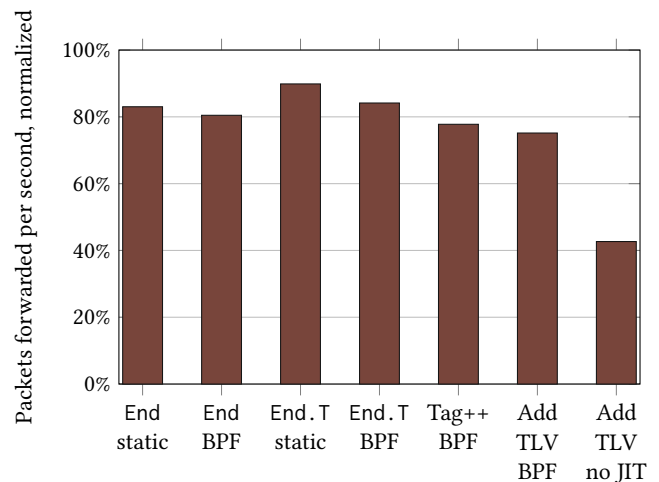


Figure 2: Simple endpoint functions are efficiently supported.

We evaluate four simple `End.BPF` functions. Our first eBPF code is `End`, i.e. a BPF function that does nothing (1 source line of code (SLOC) in its body). This function serves as a baseline to evaluate the cost of calling an eBPF function when forwarding each packet. Our second function is a BPF counterpart for `End.T`, calling `bpf_lwt_seg6_action` (4 SLOC). Our third eBPF program, `Tag++`, fetches the tag of the SRH and increments it by doing an indirect write using `bpf_lwt_seg6_store_bytes` (50 SLOC). The last one, `Add TLV`, adds an 8-byte TLV. This requires a call to `bpf_lwt_seg6_adjust_srh`, followed by a call to `bpf_lwt_seg6_store_bytes` to fill the newly allocated space with the content of the TLV (60 SLOC).

Figure 2 shows the relative forwarding performance of these functions compared to raw IPv6 forwarding. Each point is the

average over 20 measurements. The normalised standard deviation was below 1%. Compared to the static `seg6local` implementation directly written in the kernel, the eBPF equivalent of `End` has a reduced throughput of only 3%. In addition, `Tag++` decreases this throughput by 3%, by fetching and incrementing the tag through one helper call. The equivalent of `End.T` is capable of forwarding a 5% smaller throughput than its static counterpart. Finally, `Add TLV` forwards a 5% lower throughput than `End` written in BPF. `Add TLV` and `Tag++` do not have static counterparts in the `seg6local` infrastructure. In all four cases, the performance overhead is deemed acceptable.

Moreover, we use `Add TLV` to evaluate the benefits of using the eBPF JIT to support SRv6 network programming. When disabling the JIT compiler, the throughput going through `Add TLV` is divided by a factor of 1.8. Similar factors have been observed when evaluating the impact of the JIT compiler on other programs with similar complexities. This factor is expected to increase when the number of instructions per BPF program also increases.

4 USE-CASES

By leveraging `End.BPF`, network operators can implement their own in-network functions that are applied on a per-packet basis. As `End.BPF` is implemented on Linux, it can be used in a variety of environments. One important use case are the low-end access routers that are placed in many homes. We illustrate the flexibility of `End.BPF` by demonstrating three very different use cases in this section.

4.1 Passive monitoring of network delays

Delay is one of the most important performance metrics in a network. Network operators use a variety of techniques, ranging from simple ping to more precise measurements [8, 10]. While solutions such as [29, 43] have been proposed to measure the latency in datacenter networks, we propose a solution that can be deployed in the edge, up to the Customer-Premises Equipment (CPE). To illustrate `End.BPF`, we implement the recently proposed one-way delay (OWD) delay measurement for SRv6 [8] through our eBPF interface.

Our solution uses two eBPF programs installed at both tips of the path being monitored. On the router at the beginning of the path, a BPF LWT program is executed for each packet towards the given destination. This program encapsulates, using the `bpf_lwt_push_encap` helper, a defined percentage (or *probing ratio*) of the incoming regular IPv6 packets with an SRH. This SRH contains a *Delay Measurement* (DM) TLV, with a 64-bit timestamp inserted by the router that processed the packet, and a second TLV containing the IPv6 address and UDP destination port of the controller that collects the delay measurements. The segment list enforces the path on which the delay is monitored. The last segment corresponds to the router at the end of the monitored path, with the `End.DM` instruction. The SRH is built by the BPF program, the transmission timestamp is retrieved using a generic helper that we added to the Linux kernel. This function is written in 130 SLOC.

`End.DM` is an SRv6 network function implemented using `End.BPF`. At the beginning of its execution, it fetches the RX software timestamp, i.e. the time the packet left the NIC driver and entered the kernel. It subsequently inspects the SRH to retrieve the TX timestamp inside the DM TLV, as well as the TLV containing the address of the controller. Both timestamps and the information regarding the controller are sent to a user space daemon using a `perf` event since an eBPF program is not capable of sending out-of-band replies. Finally, it decapsulates the outer IPv6 header using `bpf_lwt_seg6_action` with the `End.DT6` action, and indicates that the inner IPv6 packet should be forwarded normally.

Our user space daemon is implemented in Python, it continuously listens for `perf` events. When an event is received, it creates a new thread to send both timestamps to the indicated controller in a single UDP datagram. The implementation uses the `bcc` framework [27], a BPF front-end in Python giving straightforward access to `perf` events, and is written in 100 SLOC.

We evaluated the performance impact of both BPF programs using the setup described in 3.2. R executes the `End.DM` and transit behaviour eBPF programs. S1 uses `pktgen` to generate IPv6 packets without SRH, and `trafgen` for packets with a DM TLV. The results are presented in Figure 3.

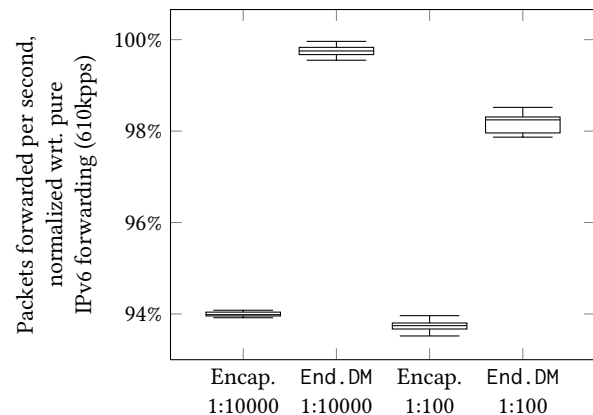


Figure 3: Impact of both BPF programs on the forwarding performances, for two probing ratios.

These results show that our passive delay monitoring is executed with almost no impact on the forwarding performances, even with a probing ratio of 1:100. We note that the transit behaviour forwards only 5% fewer packets than the native IPv6 datapath. `End.DM` has virtually no impact on performances for a 1:10000 probing ratio, even considering that all packets with a DM TLV are decapsulated.

4.2 Hybrid Access Networks

In order to offer higher bandwidths to their clients, ISPs have started to deploy hybrid access networks [19], i.e. networks that combine different access links such as xDSL and LTE. In one deployment, described in [34], a hybrid CPE router with xDSL and LTE is connected to an aggregation box with GRE Tunnels. The tunnels ensure that the packets sent by the hybrid CPE are routed to the aggregation box that reorders them. Since SRv6 inherently allows controlling

over which path each packet is forwarded, we wondered whether it is possible to leverage SRv6 BPF to design and implement a link aggregation solution that achieves good performance on home routers (CPE).

As [19], we use an aggregation box deployed in the ISP backbone. For each IPv6 route towards a client, a *LWT eBPF* program is installed on the aggregation box. It encapsulates the packets towards the client into an IPv6 header with an SRH. They are then decapsulated upon reception by the CPE and forwarded to the destination in the client’s LAN. The SRv6 decapsulation is natively performed by the kernel. The CPE uses the same eBPF program to encapsulate its packets with an SRH towards the aggregation box.

The eBPF program leverages the *LWT eBPF* hook and `bpf_lwt_push_encap` to encapsulate an SRH. Our implementation, written in 120 SLOC, performs a per-packet Weighted Round-Robin (WRR) scheduling to aggregate the bandwidths of two links. The weights of the WRR match the uplink links capacities, as seen by the CPE or the aggregation box. We use maps to store the scheduler state, i.e. the weights and the last chosen path.

To experiment with such hybrid access networks, we first configured our network as the setup 2 in Figure 1. S1, S2, A and R are the same servers as used in 3.2. S1 and S2 acts as end hosts. Node A acts as the aggregation box. Our CPE, M, is a Turrus Omnia router with a 1.6 GHz dual-core ARMv7 processor and 1 Gbps NICs. It is recent and runs OpenWRT (hence easily modifiable), making it an ideal subject for our use-case. R uses `tc netem` to insert latency on the links and to limit their bandwidth. We configure one link with a bandwidth of 50 Mbps, an average RTT of 30 ms and a standard deviation of 5 ms. The other has a bandwidth of 30 Mbps, an average RTT of 5 ms with a standard deviation of 2 ms. These values mimic current metrics of average broadband access networks.

We use our lab (Figure 3.2) to study the impact on the forwarding performances of the SRH encapsulation done in BPF and of the decapsulation performed by the kernel. UDP flows are generated between the end hosts using `iperf3` with different payload sizes at a 1 Gbps rate. The results are shown in Figure 4. The Turrus Omnia is always the bottleneck. The decapsulation induces a 10% overhead. The eBPF WRR is running without the JIT compiler, because of a bug in the current ARM32 implementation. As a consequence, the eBPF interpreter, which heavily consumes CPU resources, is the bottleneck. The setup is however almost capable of reaching the baseline performance for 1400-byte payloads, the $1.8 \times$ speedup factor provided by the JIT compiler, as demonstrated in Subsection 3.2, could be leveraged here with a functioning ARM32 implementation.

Our first experiments with TCP in this environment were a disaster. Despite an aggregated bandwidth of 80 Mbps, the TCP goodput reported by `nttcp` could only reach 3.8 Mbps. This low TCP performance is due to the difference in delays over the two links that cause TCP reordering. Commercial solutions for hybrid access networks reorder the out-of-sequence packets by using either sequence numbers in the GRE tunnels [34] or Multipath TCP proxies [13] on the CPE and the aggregation box.

Instead of using sequence numbers as in [34], we mitigate reordering by delaying the link with the lowest latency. We extend our `End.DM` implementation to handle two-way delay (TWD) measurements and deploy it on the CPE. Instead of being decapsulated by `End.DM`, the TWD probes have as last segment the IPv6 address

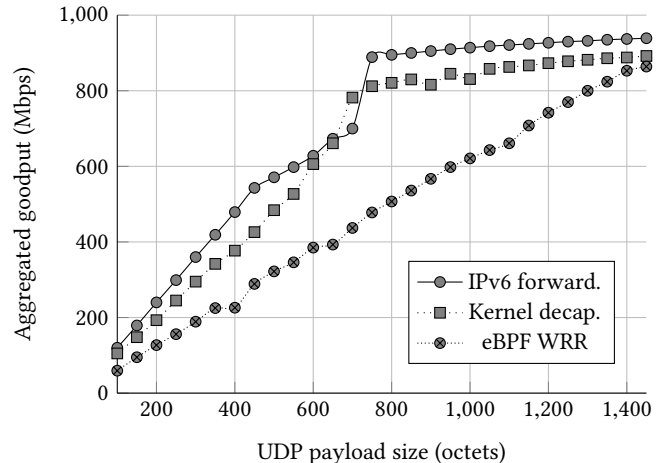


Figure 4: Aggregated UDP goodput with Turrus Omnia.

of the querier. A daemon running on the aggregation box sends TWD measurements at regular intervals to the CPE on both links. The probes return, and the daemon computes the difference of delays between the two links. Our daemon then applies a `tc netem` queuing discipline to delay the packets on the fastest path using the difference between latencies that it computed. This strategy does not fully prevent re-ordering, but still enables TCP flows to attain acceptable aggregated goodputs on links with different latencies.

We then generated TCP connections using `nttcp`. Thanks to the delay compensation, TCP could efficiently utilize the two links. A single TCP connection reached reach on average 68 Mbps while four parallel connections reached 70 Mbps.

4.3 Querying ECMP nexthops

Our third example is an enhanced version of traceroute. Given the prevalence of Equal Cost Multipath (ECMP) [30], it becomes more and more difficult for network operators to inspect routing problems with variants of traceroute [11]. Using `End.BPF`, we developed the `End.OAMP SRv6 eBPF` function which, when triggered by a packet, performs a FIB lookup to query the ECMP nexthops for its destination address and returns them to an address indicated in a TLV by the prober. Our modified traceroute leverages if possible this function at each hop, and otherwise falls back to the legacy ICMP mechanism. The eBPF function is written in 60 SLOC, whereas our custom helper returning the ECMP nexthops for a given address required only 50 SLOC in the kernel. This example underlines that, in order to extend the set of functionalities accessible to eBPF programs, new helpers can easily be added to the kernel.

5 CONCLUSION

Network programmability is high on the wish list of many network operators. In this paper, we propose, implement and evaluate an extension of the Linux implementation of IPv6 Segment Routing that enables in-network programming. We provide an eBPF interface and a set of helper functions that enable them to write their own eBPF functions and attach them to specific SRv6 segments.

Our measurements indicate that our eBPF functions have a minimal overhead compared to their static variants. Their main benefit is that they are generic. We illustrate this with three very different use cases (delay measurements, hybrid access networks and a multipath-aware traceroute). Our eBPF extensions open new ways for network operators and researchers to implement in-network functions.

Software artefacts

We release our extensions to the Linux kernel as well as the eBPF code developed for the different use cases under a GPL license: <https://github.com/Zashas/Thesis-SRV6-BPF>.

ACKNOWLEDGEMENTS

This work was partially supported by a Cisco URP grant and by CFWB within the ARC-SDN project.

REFERENCES

- [1] 2018. BPF and XDP Reference Guide. <http://cilium.readthedocs.io/en/latest/bpf/>. (2018). [Online; accessed 8 June 2018].
- [2] 2018. BPF helpers - Documentation. <https://github.com/qmonnet/bpf-helpers/blob/master/out/bpf-helpers.rst>. (2018). [Online; accessed 8 June 2018].
- [3] 2018. Linux Kernel Documentation - Linux Socket Filtering aka Berkeley Packet Filter (BPF). <https://www.kernel.org/doc/Documentation/networking/filter.txt>. (2018). [Online; accessed 8 June 2018].
- [4] 2018. Linux Weekly News - A thorough introduction to eBPF. <https://lwn.net/Articles/740157/>. (2018). [Online; accessed 8 June 2018].
- [5] 2018. The LLVM Compiler Infrastructure - Project website. <https://llvm.org/>. (2018). [Online; accessed 9 June 2018].
- [6] Ahmed AbdelSalam, Francois Clad, Clarence Filsfils, Stefano Salsano, Giuseppe Siracusano, and Luca Veltri. 2017. Implementation of virtual network function chaining through segment routing in a linux-based nfv infrastructure. In *Network Softwarization (NetSoft), 2017 IEEE Conference on*. IEEE, 1–5.
- [7] Ahmed Abdelsalam, Stefano Salsano, Francois Clad, Pablo Camarillo, and Clarence Filsfils. 2018. SERA: SEgment Routing Aware Firewall for Service Function Chaining scenarios. In *IFIP Networking 2018*.
- [8] Ali et al. 2018. *Performance Measurement in Segment Routing Networks with IPv6 Data Plane (SRv6)*. Internet-Draft draft-ali-spring-srv6-pm-02.
- [9] Zafar Ali, Clarence Filsfils, et al. 2018. *Operations, Administration, and Maintenance (OAM) in Segment Routing Networks with IPv6 Dataplane (SRv6)*. Internet-Draft draft-spring-srv6-oam-01.
- [10] Guy Almes, Sunil Kalidindi, Matthew J. Zekauskas, and Al Morton. 2016. A One-Way Delay Metric for IP Performance Metrics (IPPM). RFC 7679. (Jan. 2016). <https://doi.org/10.17487/RFC7679>
- [11] Brice Augustin, Xavier Cuvellier, Benjamin Orgogozo, Fabien Viger, Timur Friedman, Matthieu Latapy, Clémence Magnien, and Renata Teixeira. 2006. Avoiding traceroute anomalies with Paris traceroute. In *Proceedings of the 6th ACM SIGCOMM conference on Internet measurement*. ACM, 153–158.
- [12] Ahmed Bashandy, Clarence Filsfils, et al. 2018. *Topology Independent Fast Reroute using Segment Routing*. Internet-Draft draft-bashandy-rtgwg-segment-routing-tifa-05.
- [13] Olivier Bonaventure and S Seo. 2016. Multipath TCP deployments. *IETF Journal* 12, 2 (2016), 24–27.
- [14] Daniel Borkmann. 2016. On getting tc classifier fully programmable with cls bpf. *Proceedings of netdev* (2016).
- [15] Pat Bosshart, Dan Daly, Glen Gibb, Martin Izzard, Nick McKeown, Jennifer Rexford, Cole Schlesinger, Dan Talayco, Amin Vahdat, George Varghese, and David Walker. 2014. P4: Programming Protocol-independent Packet Processors. *SIGCOMM Comput. Commun. Rev.* 44, 3 (July 2014), 87–95.
- [16] Ken Calvert. 2006. Reflections on network architecture: an active networking perspective. *ACM SIGCOMM Computer Communication Review* 36, 2 (2006), 27–30.
- [17] Gaurav Dawra, John Brzozowski, John Leddy, and Clarence Filsfil. 2017. SRv6: Network as a Computer and Deployment use-cases. URL: https://pc.nanog.org/static/published/meetings/NANOG71/1445/20171005_Dawra_Segment_Routing_Ipv6_v1.pdf. (10 2017). NANOG71, San Jose, CA.
- [18] Yoann Desmoucheaux, Pierre Pfister, Jérôme Tollet, Mark Townsley, and Thomas Clausen. 2017. SRLB: The Power of Choices in Load Balancing with Segment Routing. In *ICDCS 2017*. IEEE, 2011–2016.
- [19] G. Fabregas (Ed). 2016. TR-349: Hybrid Access Broadband Network Architecture. (July 2016). Broadband Forum.
- [20] David Lebrun Fabien Duchène and Olivier Bonaventure. 2018. SRv6Pipes: enabling in-network bytestream functions. In *IFIP Networking 2018*.
- [21] David C Feldmeier, Anthony J McAuley, Jonathan M Smith, Deborah S Bakin, William S Marcus, and Thomas M Raleigh. 1998. Protocol boosters. *IEEE Journal on Selected Areas in Communications* 16, 3 (1998), 437–444.
- [22] Clarence Filsfils et al. 2018. *SRv6 Network Programming*. Internet-Draft draft-filsfils-spring-srv6-network-programming-05.
- [23] Clarence Filsfils, Francois Clad, Pablo Camarillo, Jose Liste, Prem Jonnalagadda, Milad Sharif, Stefano Salsano, and Ahmed AbdelSalam. 2017. IPv6 Segment Routing. In *SIGCOMM'17, Industrial demos*.
- [24] Clarence Filsfils, Nagendra Kumar Nainar, Carlos Pignataro, Juan Camilo Cardona, and Pierre Francois. 2015. The segment routing architecture. In *Global Communications Conference (GLOBECOM), 2015 IEEE*. IEEE, 1–6.
- [25] Clarence Filsfils, Stefano Previdi, Les Ginsberg, Bruno Decraene, Stephane Litkowski, and Rob Shakir. 2018. *Segment Routing Architecture*. RFC 8402. RFC Editor.
- [26] Thomas Graf. 2016. bpf: BPF for lightweight tunnel encapsulation. <https://lwn.net/Articles/708020/>. (2016). [Online; accessed 8 June 2018].
- [27] Thomas Graf. 2018. GitHub repository of the bcc project. <https://github.com/iovisor/bcc>. (2018). [Online; accessed 8 June 2018].
- [28] Brendan Gregg. 2018. perf Examples. <http://www.brendangregg.com/perf.html>. (2018). [Online; accessed 8 June 2018].
- [29] Chuanxiong Guo, Lihua Yuan, Dong Xiang, Yingnong Dang, Ray Huang, Dave Maltz, Zhaoyi Liu, Vin Wang, Bin Pang, Hua Chen, Zhi-Wei Lin, and Varugis Kurien. 2015. Pingmesh: A Large-Scale System for Data Center Network Latency Measurement and Analysis. *SIGCOMM Comput. Commun. Rev.* 45, 4 (Aug. 2015), 139–152.
- [30] Christian Hopps. 2000. Analysis of an equal-cost multi-path algorithm. RFC 2992. (2000).
- [31] Diego Kreutz et al. 2015. Software-defined networking: A comprehensive survey. *Proc. IEEE* 103, 1 (2015), 14–76.
- [32] David Lebrun and Olivier Bonaventure. 2017. Implementing IPv6 Segment Routing in the Linux Kernel. In *Applied Networking Research Workshop 2017*. See <https://irtf.org/anrw/2017/anrw17-final3.pdf>.
- [33] David Lebrun, Mathieu Jadin, François Clad, Clarence Filsfils, and Olivier Bonaventure. 2018. Software resolved networks: Rethinking enterprise networks with IPv6 segment routing. In *SOSR'18*. ACM, 6.
- [34] N. Leymann, C. Heidemann, M. Zhang, B. Sarikaya, and M. Cullen. 2017. *Huawei's GRE Tunnel Bonding Protocol*. RFC 8157. RFC Editor.
- [35] Steven McCanne and Van Jacobson. 1993. The BSD Packet Filter: A New Architecture for User-level Packet Capture.. In *USENIX winter*, Vol. 93.
- [36] Nick McKeown, Tom Anderson, Hari Balakrishnan, Guru Parulkar, Larry Peterson, Jennifer Rexford, Scott Shenker, and Jonathan Turner. 2008. OpenFlow: enabling innovation in campus networks. *ACM SIGCOMM Computer Communication Review* 38, 2 (2008), 69–74.
- [37] Rui Miao, Hongyi Zeng, Changhoon Kim, Jeongkeun Lee, and Minlan Yu. 2017. SilkRoad: Making Stateful Layer-4 Load Balancing Fast and Cheap Using Switching ASICs. In *SIGCOMM '17*. 15–28.
- [38] Stefano Previdi, Clarence Filsfils, et al. 2018. *IPv6 Segment Routing Header (SRH)*. Internet-Draft draft-ietf-6man-segment-routing-header-14.
- [39] Muhammad Shahbaz, Sean Choi, Ben Pfaff, Changhoon Kim, Nick Feamster, Nick McKeown, and Jennifer Rexford. 2016. PISCES: A Programmable, Protocol-Independent Software Switch. In *SIGCOMM '16*. 525–538.
- [40] Jonathan M Smith, Kenneth L Calvert, Sandra L Murphy, Hilarie K Orman, and Larry L Peterson. 1999. Activating networks: a progress report. *Computer* 32, 4 (1999), 32–41.
- [41] David L Tennenhouse, Jonathan M Smith, W David Sincoskie, David J Wetherall, and Gary J Minden. 1997. A survey of active network research. *IEEE communications Magazine* 35, 1 (1997), 80–86.
- [42] David L Tennenhouse and David J Wetherall. 1996. Towards an active network architecture. *ACM SIGCOMM Computer Communication Review* 26, 2 (1996), 5–17.
- [43] Yibo Zhu, Nanxi Kang, Jiabin Cao, Albert Greenberg, Guohan Lu, Ratul Mahajan, Dave Maltz, Lihua Yuan, Ming Zhang, Ben Y. Zhao, and Haitao Zheng. 2015. Packet-Level Telemetry in Large Datacenter Networks. In *SIGCOMM '15*. 479–491.