Fabricio Eduardo Rodríguez Cesén

# Design, Implementation and Evaluation of IPv4/IPv6 Longest Prefix Match support in Multi-Architecture Programmable Dataplanes

# Projeto, Implementação e Avaliação do Suporte de Casamento com Prefixo Mais Longo para IPv4/IPv6 em Planos de Dados Programáveis Multi-Arquitetura

Campinas

2018

Fabricio Eduardo Rodríguez Cesén

# Design, Implementation and Evaluation of IPv4/IPv6 Longest Prefix Match support in Multi-Architecture Programmable Dataplanes

# Projeto, Implementação e Avaliação do Suporte de Casamento com Prefixo Mais Longo para IPv4/IPv6 em Planos de Dados Programáveis Multi-Arquitetura

Dissertation presented to the Faculty of Electrical and Computer Engineering of the University of Campinas in partial fulfillment of the requirements for the degree of Master, in the area of Computer Engineering.

Dissertação apresentada à Faculdade de Engenharia Elétrica e de Computação da Universidade Estadual de Campinas como parte dos requisitos exigidos para a obtenção do título de Mestre em Engenharia Elétrica, na Área de Engenharia de Computação.

Supervisor: Prof. Dr. Christian Rodolfo Esteve Rothenberg

Este exemplar corresponde à versão final da tese defendida pelo aluno Fabricio Eduardo Rodríguez Cesén, e orientada pelo Prof. Dr. Christian Rodolfo Esteve Rothenberg

Campinas

2018

Ficha catalográfica
Universidade Estadual de Campinas
Biblioteca da Área de Engenharia e Arquitetura
Elizangela Aparecida dos Santos Souza - CRB 8/8098

# COMISSÃO JULGADORA - DISSERTAÇÃO DE MESTRADO

**Candidato:** Fabricio Eduardo Rodríguez Cesén RA: 163682

**Data da Defesa:** 18 de dezembro de 2018

**Título da Tese:** "Design, Implementation and Evaluation of IPv4/IPv6 Longest Prefix Match support in Multi-Architecture Programmable Dataplanes

Prof. Dr. Christian Rodolfo Esteve Rothenberg (FEEC/UNICAMP)(Presidente)

Prof. Dr. Marco Aurélio Amaral Henriques (FEEC/UNICAMP)

Prof. Dr. Fábio Luciano Verdi (/UFSCar - Campus Sorocoaba)

A ata de defesa, com as respectivas assinaturas dos membros da Comissão Julgadora, encontra-se no SIGA (Sistema de Fluxo de Dissertação/Tese) e na Secretaria de PósGraduação da Faculdade de Engenharia Elétrica e de Computação.

*To my wife, Aline, who has been a constant source of support and encouragement.*

# Acknowledgements

First and foremost, I have to thank my wife for her love and support. An special thank to my family, despite the difficulties of being far away, they always have been supporting me and encouragement to continue with my dreams.

My parents, Mariana and Vinicio, who have always supported me unconditionally and whose good examples have taught me to work hard for the things that I aspire to achieve.

I would like to sincerely thank my advisor, Professor Christian, for his guidance and support throughout this study and especially for his confidence in me.

To all my friend and colleagues, they have been an important part of this process, with their advice and support.

# Abstract

Among the New trends in dataplane programmability inside Software Defined Networking (SDN) stand out the efforts to bring multi-platform support with a high definition of the information that is processed by the dataplane pipeline. The Programming Protocol-Independent Packet Processors (P4) Domain Specific Language (DSL) is an emerging trend to express how the packets are processed by the dataplanes of a programmable network platform. In parallel, OpenDataPlane (ODP) project creates an open-source, cross-platform set of Application Programming Interfaces (APIs) designed for the networking dataplane. Multi-Architecture Compiler System for Abstract Dataplane (MACSAD) is an approach to converge P4 and ODP in a conventional compilation process, achieving portability of the dataplane applications without affecting the target performance improvements.

This thesis aims at adding IPv4/IPv6 Longest Prefix Match (LPM) support to MACSAD integrated with ODP APIs and P4 programmability delivering high-performance dataplane capabilities. The proposed LPM support for MACSAD combines the lookup algorithm and the ODP API library with MACSAD table support, to create a complete forwarding base used in the LPM process. The IPv4 implementation adapts the current ODP lookup algorithm to work with MACSAD. IPv6 lookup implementation, currently not supported by ODP, it is an extension of the IPv4 support developed using the same algorithm adapted to a 128-bit key. For the performance evaluation of the LPM support, we use a traffic generator tool Network Function Performance Analyzer (NFPA) that allows generating different types of traffic across MACSAD. Experimental results show that it is possible to reach a throughput of 10G with packets sizes of 512 Bytes and above. As a side contribution on this front is the open source BB-Gen packet crafter tool.

**Keywords**: P4; Software Defined Networking; Performance analysis; Programmable networks.

# Resumo

Dentre as novas tendências em programação de dataplane dentro de SDN (Software Defined Networking) destacam-se os esforços para prover um suporte multi-plataforma dotado de alta definição das informações que são processadas pelo pipeline do plano de dados. Como forma de mitigar tais problemas, verifica-se que a Linguagem Específica de Domínio (DSL) Programming Protocol-Independent Packet Processors (P4) desponta como uma tendência emergente para expressar como os pacotes são processados pelo plano de dados de uma plataforma de rede programável. De modo independente e em paralelo, constata-se que o projeto OpenDataPlane (ODP) cria um conjunto de plataformas abertas de Application Programming Interfaces (APIs) projetado para o plano de dados de rede. Isso posto, tem-se que o Multi-Architecture Compiler System for Abstract Dataplanes (MACSAD) surge como uma abordagem para convergir P4 e ODP em um processo de compilação convencional, arquivando a portabilidade dos aplicativos de plano de dados sem afetar as melhorias de desempenho do alvo.

Este trabalho tem como objetivo adicionar o suporte do Longest Prefix Match (LPM) do IPv4/IPv6 ao MACSAD, integrado com as APIs do ODP e à programação P4, oferecendo recursos de planejamento de dados de alto desempenho. O suporte ao LPM proposto para o MACSAD combina o algoritmo de lookup e a biblioteca da API do ODP com o suporte à tabela MACSAD, para criar uma base de encaminhamento completa usada no processo do LPM. A implementação do IPv4 adapta o atual algoritmo de lookup do ODP para trabalhar com o MACSAD. A implementação de lookup IPv6, atualmente não suportada pelo ODP, é uma extensão do suporte IPv4 que é desenvolvido usando o mesmo algoritmo adaptado a uma chave de 128 bits. Para a avaliação de desempenho do suporte ao LPM, utilizamos uma ferramenta geradora de tráfego Network Function Performance Analyzer (NFPA) que permite gerar diferentes tipos de tráfego no MACSAD. Cabe ainda destacar, como uma contribuição lateral deste trabalho, o desenvolvimento da ferramenta geradora de pacote BB-Gen, já com lançamento open source. Resultados experimentais mostram que é possível atingir um throughput de 10G com tamanhos de pacotes de 512 bytes ou superiores.

**Palavras-chaves**: P4; Rede Definida por Software; Desempenho; Redes Programáveis.

# List of Figures

# List of Tables

# Acronyms

**API** Application Programming Interface.

**ARM** Advanced RISC Machine.

**BPF** Berkeley Packet Filter.

**bps** bits per second.

**DApp** Dataplane Application.

**DPDK** Data Plane Development Kit.

**DSL** Domain Specific Language.

**DUT** Device Under Test.

**eBPF** Extended BPF.

**FIB** Forwarding Information Base.

**FPGA** Field Programmable Gate Array.

**GCC** GNU Compiler Collection.

**GRE** Generic Routing Encapsulation.

**HAL** Hardware Abstraction Library.

**HLIR** High Level IR.

**IR** Intermediate Representation.

**LLVM** Low Level Virtual Machine.

**LPM** Longest Prefix Match.

**MacS** MACSAD Switch.

**MACSAD** Multi-Architecture Compiler System for Abstract Dataplanes.

**MPLS** Multiprotocol Label Switching.

**NFPA** Network Function Performance Analyzer.

**ODP** OpenDataPlane.

**OF** OpenFlow.

**OVS** OpenvSwitch.

**P4** Programming Protocol-Independent Packet Processors.

**PI** Protocol Independence.

**PISA** Protocol Independent Switch Architecture.

**pps** packets per second.

**PRT** P4Runtime.

**SAI** Switch Abstraction Interface.

**SDK** Software Development Kit.

**SDN** Software Defined Networking.

**SoC** System on a Chip.

**T4P4S** Translator for P4 Switches.

**TCAM** Ternary Content Addressable Memory.

**TTL** Time to Live.

**VXLAN** Virtual eXtensible Local Area Network.

# Contents

# 1 Introduction

Hundreds of thousands of packets are being transmitted through the Internet every second. As a consequence, the requirements on the network are increasing exponentially, and it is necessary to evaluate the actual capabilities and how to manage this traffic efficiently. A flexible, re-designable and configurable network is needed to ensure the possibility to change the features of the dataplane.

Software Defined Networking (SDN) (KREUTZ *et al.*, 2014) is an emerging network architecture where the network splits into control and forwarding plane. This migration of control, which is tightly bounded in particular network devices, is converted into available computing devices, enabling the underlying abstracted infrastructure for applications and network services capable to treat the network as a logical or virtual entity. The first communications standard interface defined between the control and forwarding layers was OpenFlow (OF). It allows direct access to manipulation of the forwarding layer of physical and virtual network devices such as routers and switches. However, OpenFlow has some limitations, for example, it needs to know the types of the packet headers introducing difficulties in implementing new protocols and headers. Therefore, to enable programming the forwarding chip to support new protocols, a new protocol is proposed with independent programming abstractions, such Programming Protocol-Independent Packet Processors (P4) (BOSSHART et al., 2014).

P4 is an open source language for expressing how the packets are processed by the pipeline of a network forwarding element. It is based on a Match+Action forwarding model, and it works together with SDN protocols. Use of P4 makes it possible to reconfigure the packet's fields, and the programmers can change the way these ones flow on the switch. What is more, with target independence, the programmers do not have to know the specifications of the underlying hardware for being able to describe packet processing functionality. With the use of protocol independent, the switches do not need to be bound to any specific network protocols (BOSSHART et al., 2014). With SDN, new protocols program the routing devices that have emerged, as well as new data plan hardware and languages, such P4. On the other hand, ODP (OPENDATAPLANE, 2013c) provides an open-source, cross-platform set of Application Programming Interfaces (APIs) for the networking dataplane. The APIs enable developers to create dataplane codes across different targets, being neutral to different vendors and platforms.

With P4 and ODP working together, it is possible to determine and program dataplanes beyond multiple targets with a common compiler system. MACSAD (PATRA *et al.*, 2016; PATRA *et al.*, 2017) is an approach to accommodate P4 and ODP through

a common compilation process delivering portability of dataplane applications without compromising target performance improvements, translating P4-defined dataplanes into high-level ODP APIs. MACSAD has a designed compiler module that generates an Intermediate Representation (IR) for P4 applications. However, there is no previously published study in where any user communities or third-parties presented ODP IP Lookup support for IPv6, and MACSAD has limited support for IPv4 forwarding. In our work, we design, implement and evaluate the IPv4/IPv6 Longest Prefix Match (LPM) support for a multi-architecture compiler system (MACSAD) capable of supporting new Domain Specific Languages (DSLs) and network platforms. The performance of the prototype is tested in different dataplane instances for a 10G setup and now it is implemented using on ODP Layer 3 forwarding base, IPv4 LPM support for MACSAD. Using the mentioned base, we adapt the table management from ODP and combine with our forwarding support. Some functions considered to integrate the lookup tables from MACSAD to the ODP and with the LPM algorithm. Notable works regarding IP Lookup implementations are present in (RéTVáRI *et al.*, 2013; SHAHHAZ et al., 2016; PONG; TZENG, 2012).

The generation of packet traces is commonly required to carry out performance evaluations and hence it becomes a relevant task for the development of new network solutions and the evaluation of the existing ones. Our analysis found a gap among open source PCAP trace generators with simple interfaces and covering the relevant requirements for rich performance experiments. Some alternatives support a complete set of protocols and allow to generate different sets of PCAP traces but do not provide important traffic flow characteristics (e.g., packet size and address distribution). Other tools support different encapsulation protocols but not Virtual eXtensible Local Area Network (VXLAN) and Generic Routing Encapsulation (GRE) together in the same tool. Although there are many tools for packet generation, the complementary creation of table traces remains limited. Table traces are necessary in a programmable network to fill the table flow configuration of the network devices.

Our need for a single packet generator tool to meet the packet trace requirements of our research use cases to evaluate the performance and scalability of programmable dataplanes (PATRA *et al.*, 2016; PATRA *et al.*, 2017; PATRA *et al.*, 2018) and identify the limitations of existing tools, motivated our work to develop a packet generator prioritizing two essential characteristics: (*i*) simple to use, (*ii*) wide protocol support and rich customizability. To that end, as our baseline, we opted for the well known Scapy python library (BIONDI, 2008) for easier extension and packet manipulation.

## 1.1   Research Goals

The following section describes the current limitations with MACSAD IP lookup process along with the research objectives.

### 1.1.1   Problem Definition

The MACSAD implementation allows performing a basic Layer 2 forwarding. However, it limits the use-cases with the switch and the capabilities to perform a Layer 3 packet forwarding. The current limitations of IP lookup support in MACSAD are the following ones:

1. **Limited support of IPv4 lookup**. There is no IPv4 lookup algorithm implementation featuring the ODP APIs, resulting MACSAD to support only Layer 2 routing. Therefore, the implementation of an IP lookup LPM support expanding the use-cases for the project has become a necessity.

2. **No IPv6 support in ODP**. ODP helper library provides support for IPv4 lookup forwarding. However, there is not an actual implementation of IPv6 lookup provided by ODP. No short-term plans related to this implementation by the ODP team, hence, a high-reliability multi-protocol (IPv4/IPv6) lookup mechanism is imperative to perform the lookup process.

3. **Performance evaluation**. An evaluation of the Switch (MACSAD) with complete IPv4/IPv6 lookup process is indispensable. It has to maintain the complete capabilities of the network. Without affecting the throughput of the network independently of the LPM process, it has to conserve a high-reliability forwarding.

### 1.1.2   Objectives

To address the identified issues, the main objective of this work is to design, implement and evaluate the Longest Prefix Match (LPM) IPv4/IPv6 support in MACSAD. To this end, the following specific objectives are identified:

- **IPv4 LPM support for MACSAD.** Adapting the IP lookup algorithm to interact with ODP APIs. To provide a mechanism integrating an algorithm for IP lookup and ODP APIs in conjunction with MACSAD to perform the LPM task. Maintaining MACSAD architecture base and using ODP APIs allow us to create a multi-target platform.

- **IPv6 LPM support for MACSAD.** To extend the support of IPv4 lookup and make adapting and integrating the IPv6 LPM process plausible. Creating a com-

plete ODP library that supports IPv6 lookup, contributing in this way with the ODP project.

- **Controller support.** To complete the lookup process and following the SDN design, the controller has to interact with MACSAD to manage the packet traffic and table actions management.

- **Workload generation.** To have a broad and deep experimental evaluation, we need to reach packet traces with different characteristics such as specific header information, number of entries and, packet size.

- **Performance evaluation.** With the lookup process, the performance of the network cannot be strongly affected. We envision that implementing the correct algorithm, the performance of MACSAD will become optimal, exploiting the full capabilities of the network. Evaluating with unidirectional traffic and the packet configurations traces, with fixed destination and source MAC addresses, IP addresses and ports.

## 1.2 Methodology and Contributions

To achieve the objectives, the following main activities were planned: (i) Literature review, (ii) Binary tree and ODP LPM library implementation, (iii) MACSAD lookup implementation, and (iv) performance evaluation.

- **Literature review.** It compels the study and analysis of the state of the art of IPv4/IPv6 LPM algorithms to support a Multi-Architecture Programmable Dataplane. The study of ODP and P4 language to be used for the LPM implementation. The study of the architecture of MACSAD and the current project's support.

- **Binary tree and ODP implementation.** After the literature review, we implement the Binary tree and the ODP API library. As a first part of the process, the IPv4 support is developed, a second part consist of the IPv6 implementation, extending the lookup support and capabilities.

- **Lookup implementation.** With the binary tree and ODP library implemented, the next step planned consist to develop the LPM support for MACSAD and the integration with the controller. The implementation separates the IPv4 and IPv6 support as two processes being developed before the corresponding binary tree and ODP implementation.

- **Performance evaluation.** Using the Network Function Performance Analyzer (NFPA) tool, we evaluate MACSAD for different IPv4/IPv6 traces and platform configurations (e.g., CPU cores, I/O).

The contributions of this dissertation can be summarized as follows:

- Implementation and experimental evaluations of IPv4 and IPv6 Longest Prefix Match (LPM) support in MACSAD.

- Prototype implementation of a new ODP Helper library for the IPv6 lookup mechanism based on the current IPv4 solution and evaluation of the performance and scalability for diverse workloads and target platform configurations.

- The performance and scalability evaluation of MACSAD pipelines analyzing the impact of varying packet sizes and memory lookup tables, and investigating the impact of increased pipeline complexity of $P4_{14}$ and $P4_{16}$ use cases ranging from Ethernet and IPv4/v6 to VXLAN-based Data Center Gateway and Broadband Network Gateway (BNG). Combining the high-level Protocol Independence (PI) programmability of P4 with the low-level cross-platform (HW & SW) APIs brought by ODP.

- A P4 CLI-based packet crafter to generate packet flows formatted as PCAP files. Supporting different standard protocols and able to create the necessary traces for network function configuration and testing from a P4 file as an Input.

- Open source artifacts. The MACSAD code including IPv4 and IPv6 use cases were open sourced. BB-Gen tool, the PCAPs, and Traces generated are open for developers to evaluate the performance of the projects, being simple to integrate for different trace formats.

Contributions of this thesis have led to seven co-authored publications (see Appendix A).

## 1.3   Text Organization

The rest of this work contains the following topics. Background information including related work as presented in Chapter 2. Chapter 3 describes our architecture proposal to design and implement of IPv4/IPv6 LPM support for MACSAD, with the details of each process involved in its development. Chapter 4 presents the packet generator tool that will be used for the performance evaluation. Chapter 5 shows our performance evaluation in different scenarios. Finally, our conclusions and future work are presented.

# 2 Literature Review

In this chapter, we review relevant literature for our research and related works about problems and solutions that carries close resemblance to our approach.

## 2.1 Background

This section defines four main concepts in our research work: P4 expressing how the pipeline processes packets, ODP as the APIs for the networking dataplane, NFPA as open-source measurement application and finally a MACSAD.

### 2.1.1 Programming Protocol-Independent Packet Processors

Protocol Independent Switch Architecture (PISA) (MCKEOWN, 2016; GUREVICH, 2015) allows a custom definition of network protocols in a switch design approach. Top-down these devices allow us to tell them how to process packets (GUREVICH, 2015). The forwarding plane does not know any protocols until it is programmed. In Figure 1 is shown the logical mapping dataplane design for physical resources. One of the principal characteristics of P4 allows to reconfigure PISA devices in the field (and eventually on the fly) (MCKEOWN, 2016).

Figure 1 – Mapping to Physical Resources. Source (GUREVICH, 2015)

Figure 2 shows the capability of P4 to map custom protocols together with defined protocols to the physical resources.

Figure 2 – Mapping custom protocol to Physical Resources. Source (GUREVICH, 2015)

P4 is an innovation providing an abstract model suitable for programming the network dataplane (BOSSHART et al., 2014). A P4-enabled device is protocol independent. It delineates the packet headers and specifies the packet parsing and processing behaviors.

A P4 program include the flowing elements:

1. **Header definitions,** that specify the field names and widths for protocol headers (see listing 2.1).

```
1  header_type ethernet_t {
2    fields {
3        dstAddr : 48;      // Destination MAC address
4        srcAddr : 48;      // Source MAC address
5        etherType : 16;     // Ethernet type
6    }
7  }
8  header ethernet_t ethernet;
```

Listing 2.1 – P4 header definition example

2. **Metadata,** holds information about the packet that is not normally represented by the packet data (see listing 2.2).

```
1  header_type local_metadata_t {
2    fields {
3      cpu_code : 8;        // Code for packet going to CPU
4      port_type : 2;       // Type of port
5      ingress_error : 1;    // An error in ingress port check
6      was_mtagged : 1;     // Track if pkt was mtagged on ingr
7    }
8  }
```

```
9    metadata local_metadata_t local_metadata;
```

Listing 2.2 – P4 metadata definition example

3. **Registers, Meters, and Counters,** for state independent of packets (see listing 2.3 and 2.4).

```
1    counter ip_pkts_by_dest {
2      type : packets;
3      direct : ip_host_table;
4    }
```

Listing 2.3 – P4 counter definition example

```
1    meter customer_meters {
2      type : bytes;
3      instance_count : 1000;
4    }
```

Listing 2.4 – P4 meter definition example

4. **Packet parser specification,** that generates a Parsed Representation where the match+action tables operate (see listing 2.5).

```
1    parser start {
2      return parse_ethernet;
3    }
4
5    parser parse_ethernet {
6      extract(ethernet);
7      return select(latest.etherType) {
8      0x0800 : parse_ipv4;
9      default: ingress;  }
10   }
11
12   parser parse_ipv4 {
13     extract(ipv4);
14     return ingress;
15   }
```

Listing 2.5 – P4 parser example

5. **Match-action table specification,** identifying the packet and metadata fields to be read and the possible actions to execute in response (see listing 2.6).

```
1    table sendout {
2      reads {
3        standard_metadata.egress_port : exact;
4      }
5      actions {
```

```
 6        on_miss;
 7        rewrite_src_mac;
 8      }
 9      size : 512;
10   }
```

Listing 2.6 – P4 Match-action table specification example

6. **Actions,** functions that may be parameterized and that invoke one or more primitives (see listing 2.7).

```
1   action on_miss() {
2   }
3
4   action rewrite_src_mac(smac) {
5     modify_field(ethernet.srcAddr, smac);
6   }
```

Listing 2.7 – P4 actions specification example

7. **Control flow,** indicating the table execution sequence, with support for conditional branching (see listing 2.8).

```
1   control ingress {
2     apply(sendout);
3   }
```

Listing 2.8 – P4 control flow specification example

Figure 3 represents the abstract forwarding model, illustrating how the pipeline of a network forwarding element should process packets (SHAHHAZ et al., 2016). When a packet is received, the headers enter into the parser pipeline then they pass through the match+action tables flow, and finally the headers return back, and the packet is sent.



Figure 3 – P4 Abstract Forwarding Model. Source: Adapted from (MCKEOWN, 2016)

In the P4 program, all these activities are defined inside the ingress and egress of control flows. This is possible due to the protocol independent feature that details custom headers expressing arbitrary network protocol headers and fields.

## 2.1.2 OpenDataPlane

ODP project is a networking dataplane API specification. It allows application developers to write and implement dataplane applications. It can leverage portability and multi-platform support, besides the use of specific hardware acceleration capabilities. ODP defines a set of high-level common APIs (e.g., CPU control, operations on a packet Input/Output interface, operations on memory), bringing the opportunity of span standard features across multi-targets (e.g. Advanced RISC Machine (ARM)v7, ARMv8, MIPS64, Power, and x86) making dataplane applications portable. Table 1 presents the ODP supported platforms.

The principal attributes of the ODP API are summarized below:

- Open Source and open for contributions.

- Vendor and platform independent.

- Application-centric, encloses the functional needs of dapaplane applications.

- Allows portability.

- Defined in conjunction with application writers and platform implementers.

- Architecture efficiently used on a wide range of different platforms.

- Sponsored, governed, and maintained by the Linaro Networking Group (LNG).

Figure 4 presents the ODP software stack. ODP is composed by a common layer (ODP APP) and the implementation layer (ODP Implementation). The Applications defined in the common layer are portable across all ODP implementations. An application to be executed by the ODP API will be linked to the ODP implementation for the specific execution platform. The purpose of the implementation layer is to provide a mapping of ODP APIs to the underlying capabilities (i.e., hardware co-processing and acceleration

Table 1 – ODP supported platforms

| Company | Supported Platforms |
|---|---|
| Cavium | ThunderX CN88xx 24-48core ARMv8<br>OCTEON TX CN83/81xx 1-24 core ARMv8 |
| Kalray | MPPA |
| Freescale | QorIQ – ARM based DPAA2 architecture LS2080, LS2085<br>QorIQ – ARM & PowerPC based DPAA architecture LS1043 |
| Texas Instrument | Keystone2 Cortex A15 |
| Linaro | PCIe NIC optimized implementation (odp-dpdk) |
| Marvell | Marvell ARMADA SoC Implementation |

Figure 4 – ODP software stack in a Linux-based target. Source (OPENDATAPLANE, 2013c)

support) of System on a Chips (SoCs) hosting ODP implementations. The implementation has been optimized for a particular hardware platform (SoC or Server). It will call the Vendor Specific Hardware Blocks and Software Libraries (Software Development Kit (SDK)) through an inline call and takes advantage of the vendor extension that it is not yet standardized (OPENDATAPLANE, 2013c).

ODP consists of a common layer and an implementation layer. Applications written to the common layer are portable across all ODP implementations. To compile and run an ODP application, it is compiled against a specific ODP implementation layer.

The purpose of the implementation layer is to provide an optimal mapping of ODP APIs to the underlying capabilities (including hardware co-processing and acceleration support) of SoCs hosting ODP implementations. As a bootstrapping mechanism for applications, as well as to provide a model for ODP implementers, ODP provides a 'linux-generic' reference implementation designed to run on any SoC which has a Linux kernel. While linux-generic is not a performance target, it does provide a starting point for ODP implementers and application programmers alike. As a pure software implementation of ODP, linux-generic is designed to provide best-in-class performance for general Linux data plane support.

ODP helper library (OPENDATAPLANE, 2013a) offers a basic support for table management (lookup mechanism), contrary to Data Plane Development Kit (DPDK), which has various fully optimized table management libraries. It also offers support to cuckoo, hash, IP lookup and linear tables.

ODP has been taken up by vendors that provide much more functionalities in their hardware than the plain NIC capabilities. For instance the ODP-DPDK (OPENDATAPLANE, 2013b) packet, as a NIC + ODP library, helping the migration from the lower

Figure 5 – Linux Kernel without DPDK vs Linux Kernel with DPDK

level of DPDK API to the ODP abstraction.

To improve the performance, ODP can use the user-space fast packet processing I/O support for traditional NICs. Odp-linux (a functional reference targeting simplicity over performance) has PKTIO support for Netmap (RIZZO, 2012) and DPDK (INTEL, 2014). With different characteristics and architecture, the DPDK and Netmap drivers are described below:

- DPDK: it is a Linux Foundation project consisting of libraries and drivers for developers to create fast packet processing Dataplane Application (DApp). It started with x86 architecture and later expanded to other platforms like ARM and IBM POWER chips, and so forth. Figure 5 compares packet path in a vanilla linux-kernel and with DPDK driver and shows that a DApp can access network interface from user space for faster packet processing.

- Netmap: its framework allows fast packet access from the network card. It reduces packet processing cost by preallocation of resources, batch packet handling and use of shared memory to achieve higher performance. It works along with Linux kernel and allows the use of Linux tools like "ethtool", "ip" for interface management.

### 2.1.3 Network Function Performance Analyzer

The NFPA (CSIKOR *et al.*, 2015a; CSIKOR *et al.*, 2015b), was proposed as a benchmarking tool that allows the user to measure important performance metrics of a network function compiled on any hardware and software combination, and to compare the results collected in a public Database. All the experiments running through the NFPA tool are following standardized methodologies (BRADNER; MCQUAID, 1999).

NFPA was implemented using Python. A configuration file can configure the measurements parameters and traffic traces to be used. NFPA is built on Intel's DPDK (INTEL, 2014) to avoid the limitation of kernel space with network card drivers; for sending and receiving traffic, NFPA uses PktGen (WILES, 2010; TURULL *et al.*, 2016) with custom Lua scripts for parameterizing, automating and controlling the measurements. The results are saved in a local SQLite database and have generated Gnuplot (GNUPLOT, 1986) graphs from the performance results.

NFPA provides a complete selection of synthetic traffic traces with several packet headers and sizes to obtain realistic scenarios as possible. There are also supported custom user tracers in Pcap files. The traffic traces available by the NFPA project (CSIKOR *et al.*, 2015a) includes Layer 2, Layer 3 (IPv4), VxLAN (MAHALINGAM *et al.*, 2014) and Multiprotocol Label Switching (MPLS) (ROSEN *et al.*, 2001).

The NFPA standalone node is connected to the Device Under Test (DUT) as it is presented in Figure 6[1]. The user can set the measurement setup (e.g., details of the hardware and software components, the number of repeated measurements and their duration) and the traffic traces to use. NFPA sends packets on port 0, they are processed and forwarded in the DUT node and are received on port 1. Then the throughput of the DUT is calculated in terms of packets per second (PPS) and bits per second (BPS). Once the measurement is finished NFPA saves the performance results in the local database.

When is evaluated the performance capabilities of an Ethernet device, the principal indicator is the raw bandwidth (BPS). However, it is also important to analyze the DUT capability to switch/route as many packets as required to achieve wire rate performance. This is the PPS metric (JUNIPER NETWORKS, 2015).

To calculate the amount of PPS to archive wire-rate performance it is necessary to consider the fact that the IP protocol allows variable payload sizes. With smaller packet sizes passing over the link, it is required more packets to achieve the wire rate performance, while if the packet size increases the necessary packets to saturate the link decreases:

$$pps = \frac{Gbps}{(pkt\_size + 20) * 8}$$

During the calculation of the PPS value, we need to consider the space that each packet will occupy, for this end, we will use the frame size, as well as the Inter-frame Gap (12 Bytes), and the Preamble (8 Bytes):

$$20 = \underset{12\ Bytes}{Interframe\ Gap} + \underset{8\ Bytes}{Preamble}$$

---

[1] For illustration purpose MACSAD is the DUT

Figure 6 – NFPA Architecture. Source: Adapted from (CSIKOR *et al.*, 2015b)

If we consider the smaller packet size (64 Bytes) and a line rate throughput of 10 Gbps. It is possible to calculate the necessary PPS to saturate the link. In this case the maximum PPS is 14,880,952:

$$pps = \frac{10\ Gbps}{(64 + 20) * 8} = 14,880,952\ pps = 14.88\ Mpps$$

In table 2 are summarized the pps line rate values for different packet sizes.

Table 2 – pps Calculation for different packet sizes

| Speed | bits/second | bytes/second | Packet Size | Maximum PPS |
|--------|----------------|---------------|-------------|-------------|
| 10 Gbps | 10,000,000,000 | 1,250,000,000 | 64 | 14,880,952 |
| 10 Gbps | 10,000,000,000 | 1,250,000,000 | 128 | 8,445,946 |
| 10 Gbps | 10,000,000,000 | 1,250,000,000 | 256 | 4,528,986 |
| 10 Gbps | 10,000,000,000 | 1,250,000,000 | 512 | 2,349,624 |
| 10 Gbps | 10,000,000,000 | 1,250,000,000 | 1024 | 1,197,318 |
| 10 Gbps | 10,000,000,000 | 1,250,000,000 | 1280 | 961,538 |
| 10 Gbps | 10,000,000,000 | 1,250,000,000 | 1518 | 812,744 |

## 2.1.4  Multi-Architecture Compiler System for Abstract Dataplanes

MACSAD aims at taking the advantages of P4 and ODP in a common compiler system to determine and program dataplanes beyond multiple targets. MACSAD can

Figure 7 – MACSAD Architecture. Source (PATRA *et al.*, 2016)



Figure 8 – Compilation Process. Source: Adapted from (PATRA *et al.*, 2016)

deliver portability of dataplane applications without compromising target performance improvements. The architecture overview is shown in Figure 7. It has three main modules: (i) Auxiliary Frontend, (ii) Auxiliary Backend and (iii) Core Compiler.

In the way that MACSAD may support different targets optimally which is a necessary solution to:

1. Support different DSLs with a plug-in framework (Auxiliary Frontend, with P4 as the premier choice).

2. Multi-platform support, linking together target-specific SDKs (Auxiliary Backend, with ODP as the desired backend).

3. Turns the intermediate representation for P4 applications generated by the auxiliary frontend, and this together with the auxiliary backend, into the imaged target (Core Compiler, compounded of a Transpiler and a Compiler).

### 2.1.4.1 Auxiliary Frontend

The Auxiliary Frontend creates the IR for the Core Compiler, based on the P4 code as an input. The p4-hlir project[2], which is supported by the P4 organization, it is used to translate the P4 programs into a High Level IR (HLIR). The Auxiliary Frontend is designed as a plug-in framework and allows to aggregate several DSLs. Figure 8 shows the generation of the HLIR from the P4, using the p4-hlir support.

### 2.1.4.2 Auxiliary Backend

To give a common SDK for the compiler, it is used the Auxiliary Backend. The compiler incorporates the ODP APIs. To create the connection between P4 and ODP abstraction, the the development libraries are necessary. These libraries allow the packet processing inside the dataplane pipeline. This module provides the API auto-generation support, and enables control protocols like Switch Abstraction Interface (SAI)[3], Open-Flow[4], etc.

### 2.1.4.3 Core Compiler

As the principal part of MACSAD Switch (MacS), the Core Compiler encompasses the Transpiler and Compiler internal modules. With the HLIR generated by the Auxiliary Frontend and the ODP APIs provided by the Auxiliary Backend, it is compiled the MacS. Figure 8 displays the complete Compilation Process.

The **Transpiler** takes the input from the Auxilary Frontend (HLIR) and auto-generates the Datapath Logic codes. In Figure 8 is presented the Transpiler's processing. The Datapath Logic is defined using 'C' language, it is required by the pipeline and used by the Compiler. The Transpiler determines the lookup mechanism, the size, and type of tables that are going to be created, with the resources available on the target. The 'C' codes generated have information from the P4 program, ODP APIs (SDKs), datapath definitions, and helper libraries. The Transpiler allows enabling some code optimizations for example 'Dead Code Elimination' identifying reachability in a dependency graph of parser logic created.

The **Compiler** creates the MacS for the desired target with the 'C' codes generated. With MACSAD is supported the GNU Compiler Collection (GCC) and Low Level Virtual Machine (LLVM) based compiler, supporting multiple targets and optimization tools.

---

[2]  https://github.com/p4lang/p4-hlir
[3]  https://github.com/opencomputeproject/SAI
[4]  https://www.opennetworking.org/sdn-resources/openflowI

## 2.1.5   Helper APIs

MACSAD is implemented with the use of a number of APIs (Auxiliary Backend module). These APIs helps the Compiler submodule in the compilation process, covering the Parser and Table functionalities of a P4 program over the ODP SDKs. An example of this Helper APIs or libraries that MACSAD can use is the IPv4 LPM implementation from ODP, adding new features and use cases to the project.

## 2.2   Related Work

In networking devices, IP lookup forms a bottleneck during the packet forwarding due to the lookup speed unable to deal with the increase in link bandwidth (RAVIKU-MAR; MAHAPATRA, 2004).

In PISCES (SHAHHAZ et al., 2016), the packet forwarding and processing are presented in a high-level DSL as P4, compiling down to run on the underlying software switch. The implementation is not protocol specific, which means new features can be added. The researchers compared the performance to projects as OpenvSwitch (OVS) (OVS, 2009) a hardwired hypervisor switch. PISCES program might be optimized about 40 times compared as an equivalent program in Open vSwitch.

The work in (RéTVáRI et al., 2013) shows how to write the IP Forwarding Information Base (FIB) to make forwarding decisions, with basically zero cost on LPM and FIB update. They extend a static entropy compressed FIB representation (based on the labeled tree entropy measure (FERRAGINA et al., 2009)) with optimal lookup. The authors evaluate the algorithm on a Linux kernel prototype and an FPGA implementation. The compressors encode a FIB of more than 440K prefixes to about 100–400 Kbytes of memory, with an increase in lookup throughput and without time reductions with FIB updates. Later, the researchers re-designed the prefix tree, used commonly for IP lookup, to support and update lookups in optimized time. They compare the work with a Trie-based FIB schemes kernel implementation of Patricia trees (SKLOWER, 2001), the trie-based consumes 24 bytes per node and a single IP lookup cost 32 random memory accesses.

Related with the Compressing IP Forwarding Tables: Towards Entropy Bounds and Beyond work, CUCKOOSWITCH (ZHOU et al., 2013) presents a Scalable High Performance Ethernet Forwarding. The project is an Ethernet switch, with a design based on software, and created around memory efficiency, high-performance, and highly-concurrent hash table for compact and fast FIB lookup. The presented switch can process 92.22 million packets (minimum sized) per second, configured with eight 10 Gbps Ethernet interfaces on a commodity server. The performance test is also run with a continues maintenance of the forwarding table with one billion entries. The evaluation process is

performed in three main steps: (i) they evaluate the forwarding with no switching or FIB lookup involved, these results as a base line for the capacity of the platform using DPDK, (ii) they evaluate the proposed optimizations contribution to the performance of the hash tables and the full system forwarding packets, (iii) the project is compared with other common hash table implementations (INTEL, 2011) (hash table micro-benchmarks and complete system evaluation).

Related to IPv4 and IPv6 forwarding implementations, the work (NIKOLENKO *et al.*, 2016) proposes an abstraction layer able to represent IPv6 FIBs on existing IP and even MPLS infrastructure. Due to most of the forwarding methods that efficiently represent IP-based FIBs do not scale well to IPv6 (larger 128-bit address width) (EATHERTON GEORGE VARGHESE, 2004; RéTVáRI *et al.*, 2013; SRINIVASAN; VARGHESE, 1999). The implementations are common in form of decision tree representation for IPv6 software FIB. The issue with these implementations is that prefix trees are inherently sequential. For this reason, LPM lookup involves multiple consecutive steps, being the total number of steps not optimal for IPv6 (128-bit address). Each step includes separated random access to memory, limiting the total number of levels that the tree can take. For this reason, most of IPv4 FIB trees implementations are not extensible to IPv6 with optimal time and memory requirements. To avoid this trouble, the project appeals to an IPv6 FIB representation on a parallel implementation, lookup classified prefixes into groups, where each group contains all prefixes of the same length and perform the lookup of each group simultaneously.

The work (PONG; TZENG, 2012) shows a distinct LPM lookup scheme to reach concise lookup tables. The project is based on unified hash tables and can handle IPv4 and IPv6 lookup simultaneously. They simplify the table format to earn a better prefix aggregation, also being simplified the implementation process. Due to the hash table implementation, multiple possible buckets are permitted, allowing parallel search over tables during the LPM lookup process. Each lookup takes two cycles on average to complete a lookup and can handle 250 millions of packets per second. A related work (TOBOLA; KOřENEK, 2011) implements a Hash-Tree Bitmap algorithm for fast LPM, also supporting IPv4 and IPv6 lookups. The fast hash implementation allows passing through a limited part of the IP prefix tree. The proposed algorithm uses 16 parallel hash units and two consequent tree Bitmap lookups. The throughput archived by the algorithm reaches 100 Gbps on Virtex 5 FPGA[5] without memory and resources limitations.

For the lookup process, it is commonly used Ternary Content Addressable Memory (TCAM) to facilitate fast IP lookup as it is presented in (HANNA *et al.*, 2011) TCAMs are power-eager, expensive, and not scalable. They perform a forwarding table in tries data structures that are acceded by pipeline. The algorithm proposed is based on a multi-

---

[5]   https://www.xilinx.com/support/documentation/data_sheets/ds100.pdf

bit tree architecture and can reach a throughput of 3.1 Terabits per second. The algorithm archives a better compression ratio, compared with other proposed algorithms (WANG; TZENG, 2006; DEGERMARK *et al.*, 1997; EATHERTONK *et al.*, 2004).

(FIESSLER *et al.*, 2017) presents a Hybrid On-chip Matching combining a highly optimized configuration specialized and thus energy and resource efficient classification circuit with the capability of quickly updated and network packet process at link rate on a Field Programmable Gate Array (FPGA). The evaluation demonstrates that the hybrid implementation benefices the process, resulting in a matching engine that can manage the updates efficiently with a lower hardware resource and power dissipation footprint. The algorithm presented in this work is not restricted to a particular use-case. It can be adapted to an environment where two matching engines with different matching capabilities need to be used.

The project Translator for P4 Switches (T4P4S) (LAKI *et al.*, 2016) is a high speed packet forwarding compiled from protocol independent dataplane specifications, it uses as the auxiliary backend DPDK contrary to MACSAD that uses ODP. T4P4S runs the target independent P4 core in top of a Hardware Abstraction Library (HAL) to improves portability. The developed use cases covered L2 and L3 (IPv4) switch examples.

Primary related projects around SDN, dataplane and IP lookup implementations with IPv4/IPv6 support mentioned above are summarized in Table 3.

Here, we summarize a selected set of related PCAP trace generator solutions. Each tool has their advantages and disadvantages depending on their architecture, including features or supported platforms.

RWS (KNUTSSON, 2014) PCAP generator is based on a simple packet descriptor language. The user defines the header fields for the packets required and feed it to RWS to generate the PCAP. It can also generate invalid packets which is rather uncommon among PCAP generators. An example of an invalid packet can be a TCP packet tunneled inside a Teredo tunnel and sent over GTP-u.

Ostinato (OSTINATO, 2010) is one of the most powerful packet crafter, network traffic generator and analyzer with complete GUI support. It implements most of the common standard protocols to facilitate traffic generation and analysis. With a complex user interface and numerous feature combinations, Ostinato presents a steep learning curve to tackle with, which makes it difficult and time-consuming for users to understand and take advantage of the tool.

Scapy (BIONDI, 2008) is a packet manipulation program with Python interpreter disguised as a DSL. It can create and decode packets of an extended number of protocols. It can send and capture network traffic too. Its extended features also include some basic network tasks (e.g., scanning, trace routing, probing, arpspoof, arp-sk, arping, tcpdump,

and tethereal.). Scapy can stand out among competitors with its unique ability to arrange protocol headers in a custom sequence which may not confirm to any protocol logic. This feature allows Scapy to create invalid frames by combining techniques (e.g., VOIP decoding on WEP encrypted channel, and so on.), similar to RWS. Among other features, Scapy allows to set values for all header fields, payload, and padding. Moreover, it allows writing a list of packets to a PCAP file.

With the increasing use of P4 to define the dataplane structure and the efforts to develop complex use cases and applications, the tools to generate traffic, use cases and analyze the code are emerging. Some of the main works for the generations of traces and evaluation are summarized below.

A tool to validate a P4 program by generating test input packets and tables entries is P4pktgen (NÖTZLI *et al.*, 2018). This tool automatically generates test cases for P4 programs using symbolic execution, allowing in parallel to detect bugs in the files. A similar project is the P4app (P4.ORG, 2013). It allows to perform functional tests for a P4 program using BMV2 simple_switch in Mininet environment, being limited to a simulated environment. For a complex test the work P4 benchmark (DANG *et al.*, 2017) test the target compilers by the generation of different P4 applications with variable complexities.

## 2.3   Summary

This chapter summarized the background of main concepts that sets up the context of the dissertation and the related work and projects that are aligned with this dissertation. This dissertation's approach adds the IPv4 and IPv6 LPM support to MACSAD, creating a complete compiler system with a P4 dataplane base, being different to other approaches that lack the support of IPv6 or do not have a compiler system and a P4 base.

We briefly introduced four main concepts that form an essential part in the dissertation. P4 defining the packets pipeline process, ODP as the APIs and libraries for the dataplane, NFPA as a benchmarking tool, and finally MACSAD to join together P4 and ODP in a common compilation process, creating the software switch image.

Table 3 – Feature comparison list of different IP lookup projects

| Project | IPv4 | IPv6 | Compiler | P4 Dataplane | Target | Remarks |
|---|---|---|---|---|---|---|
| PISCES | Yes | No | Yes | Yes | Software Switch | OVS Based |
| Compressing IP Forwarding Tables: Towards Entropy Bounds and Beyond | Yes | No | No | No | General-Purpose Processor/ FPGA | Based on prefix tree |
| CUCKOOSWITCH | Yes | No | Yes | No | General-Purpose Processor | Based on cuckoo hashing |
| IPv6 Forwarding Tables on IPv4 or MPLS Dataplanes | Yes | Yes | No | No | Limited by DPDK | Parallel prefix trees implementation |
| Effective Hash-based IPv6 Longest Prefix Match | Yes | Yes | No | No | General-Purpose Processor/ FPGA | Hash-Tree Bitmap architecture |
| Concise Lookup Tables for IPv4 and IPv6 Longest Prefix Matching in Scalable Routers | Yes | Yes | No | No | Multi-Target | Based on unified hash tables |
| A Novel Scalable IPv6 Lookup Scheme Using Compressed Pipelined Tries | Yes | Yes | No | No | Multi-Target | Based on trie data structures |
| T4P4S | Yes | No | Yes | Yes | Multi-Target | X86 & ARMv8 support available |
| MACSAD (before IPv4/IPv6 Support) | No | No | Yes | Yes | Multi-Target | X86 & ARMv8 support available |
| MACSAD (after or proposed IPv4/IPv6 support) | Yes | Yes | Yes | Yes | Multi-Target | X86 & ARMv8 support available |

# 3 Design and Implementation of IPv4/IPv6 Longest Prefix Match support

Based on the problems identified in chapter 1 of **Limited support of IPv4 lookup** and **No IPv6 support in ODP**, and in the objectives of **IPv4 and IPv6 LPM support for MACSAD** including the **Controller support**, in the following section, we will describe the implementation choices and the processes embedded in MACSAD IPv4/IPv6 LPM support.

Firstly, we present the IPv4 LPM support implementation with a brief view of the MACSAD and ODP library functions. Secondly, we detail the IPv6 implementation with a complete description of the functions and processes, considering that the base of the application is the same of IPv4 but adding support for a big key size (128 bytes). Thirdly, it is depicted the Controller and following the P4 Code. Finally, we illustrate the LPM diagram flow, summarizing the MACSAD and ODP functions.

## 3.1 IPv4 LPM support

ODP helper library has implemented an IP lookup algorithm to perform the LPM process. We use the library as support for MACSAD IPv4 lookup implementation. From MacS table control it is called the function that performs the lookup process.

The IP lookup algorithm uses a binary tree to detect the overlap prefix (matching table entries). In a binary prefix tree data structure, each node has at the most two children referred as the left child and the right child. Figure 9 presents the Binary prefix tree description including the relationship through the different tree levels. Splitting the process into different levels and limiting the number of tables, we can reduce memory consumption and maintain a good lookup speed. For the ODP lookup implementation, the number of level 1 entries is 16, and the size of one level 2 and level 3 subtrees are 8. For our implementation we maintain the ODP library definition of levels, fixing the root node in 16, covering the most common prefixes length distributions (RIPE NCC, 2010), and adding two additional levels to complete the 32 bits prefix size. The number of level 2 and level 3 entries (subtrees) per cache cube is 13, and the number of prefix tree nodes per cache cube is 20.

The lookup algorithm uses a structure that stores an entry of the IP prefix table. Because of the leaf pushing, each entry of the table must have either a child entry or a nexthop information. If the child is 0 and the index is valid, this entry has a nexthop

Figure 9 – Binary prefix tree levels. Source: Adapted from (WATERLOO, 2018)

information; the index indicates the buffer that stores the nexthop value. If the child is 1, this entry has a subtree; the index indicates the buffer that stores the subtree.

The LPM lookup has three main processes, the Binary prefix tree, the ODP IP prefix lookup table, and the MACSAD lookup. To implement the lookup process, some basic methods are required depending on the process involved. The methods that we implement are mentioned in the next section along with the relationship between the different lookup process. The complete ODP code implemented can be found in Appendix B.

## 3.1.1   Binary tree

The process that involves the **Binary prefix tree** is presented in Figure 10[1]. This process is in charge of creating the prefix tree and managing the information that is going to be stored in the specific nodes. The functions involved in the process are:

- ***trie__init***, to initialize the root node of the prefix tree.

- ***trie__destroy***, to destroy the whole prefix tree (recursively).

- ***trie__insert__node***, to insert a new prefix node into the prefix tree. If the node already exists, it is updated the nexthop information, if the node does not exist the target node is created and all nodes along the path from the root to the target node.

- ***trie__delete__node***, to delete a node.

- ***trie__detect__overlap***, to detect the longest overlapping prefix.

---

[1]   For the functions name, it was maintained the ODP format using the term "trie" for the rest of the work it was used the term "tree".

Figure 10 – Binary tree

### 3.1.2   IP prefix lookup table

To integrate the Binary prefix tree process with the lookup library, the functions presented in Figure 11 are the implemented ones. This process uses the Binary prefix tree methods to originate the tables that are going to store the lookup information. The functions involved are:

- ***odph_iplookup_table_create***, to create a table.

- ***odph_iplookup_table_lookup***, to perform the table lookup.

- ***odph_iplookup_table_destroy***, to destroy the table.

- ***odph_iplookup_table_put_value***, to add a new value into the table.

- ***odph_iplookup_table_get_value***, to get a value stored into the table.

- ***odph_iplookup_table_remove_value***, to remove a entry from the table.

### 3.1.3   MACSAD lookup

MACSAD uses the IP lookup library, and the functions described above are called. The **MACSAD lookup** process is portrayed in Figure 12. This process sends the data

Figure 11 – ODP IP prefix lookup table

to be stored in the tables to the IP prefix lookup table applying the methods and the information described in the section above. The functions involved are:

- **table_create**, to create and initialize the table.

- **lpm_add**, to add a LPM value.

- **lpm_lookup**, to perform the lookup process.

- **odpc_lookup_tbls_des**, to destroy and remove all the table values.

Figure 13 shows the relationship between the MACSAD lookup functions, the ODP IP prefix lookup table and the Binary prefix tree.

## 3.2  IPv6 LPM support

### 3.2.1  Binary tree

For the IPv6 case, we maintain the logic of the ODP IPv4 **Binary prefix tree** implementation, with the additional support for 128 bits addresses. In the following list,

Figure 12 – MACSAD lookup functions

we describe the required parameters for each **Binary prefix tree** functions[2].

- ***trie_init*** function initializes the root node of the prefix tree. With *cache_get_buffer* we get a new buffer from a cache list, if there is no available buffer a new pool will be allocated. In this case, the new buffer will insert into the queue and according to the type of cache will set the initial value of the buffer.

```
1   trie_init(odph_iplookup_table_impl *tbl){
2       trie_node_t *root = NULL;
3       odp_buffer_t buffer = cache_get_buffer(tbl, CACHE_TYPE_TRIE);
4
5       if (buffer != ODP_BUFFER_INVALID) {
6           root = (trie_node_t *)odp_buffer_addr(buffer);
7           root->cidr = 0;
8           tbl->trie = root;
9           return 0;}
10      return -1;
11  }
```

Listing 3.1 – Binary tree initialization function

---

[2]   In the Listings, it was highlighted in red color the modifications between the existing ODP IPv4 implementation and the new IPv6 implementation.

Figure 13 – Lookup relationship between Binary prefix tree, ODP and MACS

- **_trie_destroy_** function can destroy the whole prefix tree (recursively), verifying the child nodes (left and right).

```
trie_destroy(odph_iplookup_table_impl *tbl, trie_node_t *trie){
  if (trie->left != NULL)
    trie_destroy(tbl, trie->left);
  if (trie->right != NULL)
    trie_destroy(tbl, trie->right);

  odp_queue_enq(
      tbl->free_slots[CACHE_TYPE_TRIE],
      odp_buffer_to_event(trie->buffer));
}
```

Listing 3.2 – Binary tree destroy function

- **_trie_insert_node_** function inserts a new prefix node into the prefix tree. If the node already exists, it is updated the nexthop information, set the return to 0 and the nexthop pointer to INVALID. If the node does not exist, the target node is created and all nodes along the path from the root to the target node, then it is set the return to 0, and the nexthop pointer points to the new buffer. If any error occurs during the process, the return will be set to -1.

```
trie_insert_node(odph_iplookup_table_impl *tbl, trie_node_t *root,
        _uint128_t ip, uint8_t cidr, odp_buffer_t nexthop)
{
  uint8_t level = 0, child;
  odp_buffer_t buf;
```

```
6     trie_node_t *node = root, *prev = root;
7     for (level = 1; level <= cidr; level++) {
8       child = WHICH_CHILD(ip, level);
9       node = child == 0 ? prev->left : prev->right;
10      if (node == NULL) {
11        buf = cache_get_buffer(tbl, CACHE_TYPE_TRIE);
12        if (buf == ODP_BUFFER_INVALID)
13          return -1;
14        node = (trie_node_t *)odp_buffer_addr(buf);
15        node->cidr = level;
16        node->parent = prev;
17        if (child == 0)
18          prev->left = node;
19        else
20          prev->right = node;
21      }
22      prev = node;
23    }
24    node->nexthop = nexthop;
25    return 0;
26  }
```

Listing 3.3 – Binary tree insert node function

- **trie_delete_node** function can be used to delete a node. However, the default prefix (root node) can not be deleted. During the process, while finding the target node all redundant nodes are removed along the path.

```
1   trie_delete_node(odph_iplookup_table_impl *tbl,
2           trie_node_t *root, _uint128_t ip, uint8_t cidr){
3     if (root == NULL)
4       return -1;
5     if (cidr == 0)
6       return -1;
7     trie_node_t *node = root, *prev = NULL;
8     uint8_t level = 1, child = 0;
9     odp_buffer_t tmp;
10    for (level = 1; level <= cidr; level++) {
11      child = WHICH_CHILD(ip, level);
12      node = (child == 0) ? node->left : node->right;
13      if (node == NULL) {
14        return -1;
15      }
16    }
17    node->nexthop = ODP_BUFFER_INVALID;
18    for (level = cidr; level > 0; level--) {
19      if (
20        node->left != NULL || node->right != NULL ||
```

```
21            node->nexthop != ODP_BUFFER_INVALID)
22          break;
23       child = WHICH_CHILD(ip, level);
24       prev = node->parent;
25       tmp = node->buffer;
26       cache_init_buffer(
27            tmp, CACHE_TYPE_TRIE, sizeof(trie_node_t));
28       odp_queue_enq(
29            tbl->free_slots[CACHE_TYPE_TRIE],
30            odp_buffer_to_event(tmp));
31       if (child == 0)
32         prev->left = NULL;
33       else
34         prev->right = NULL;
35       node = prev;
36     }
37     return 0;
38 }
```

Listing 3.4 – Binary tree delete node function

- ***trie_delete_overlap*** function detects the longest overlapping prefix. It is used to remove values from the IPv6 lookup table.

```
1  trie_detect_overlap(trie_node_t *trie, _uint128_t ip, uint8_t cidr,
2           uint8_t leaf_push, uint8_t *over_cidr,
3           odp_buffer_t *over_nexthop){
4    uint8_t child = 0;
5    uint32_t level, limit = cidr > leaf_push ? leaf_push + 1 : cidr;
6    trie_node_t *node = trie, *longest = trie;
7    for (level = 1; level < limit; level++) {
8      child = WHICH_CHILD(ip, level);
9      node = (child == 0) ? node->left : node->right;
10     if (node->nexthop != ODP_BUFFER_INVALID)
11        longest = node;
12   }
13   *over_cidr = longest->cidr;
14   *over_nexthop = longest->nexthop;
15   return 0;
16 }
```

Listing 3.5 – Binary tree detect overlap function

### 3.2.2 IP prefix lookup table

To integrate the Binary prefix tree process with the lookup library, the functions presented in Figure 11 are adapted for the IPv6 support. This process uses the Binary

prefix tree methods to originate the tables that are going to store the lookup information. The processes involved are:

- ***odph__iplookupv6__table__create*** function creates the lookup table. This function, firstly, is going to check if the parameters are valid. Then it is going to guarantee that there are no existing tables with the same name using the *odph_iplookup_table_lookup* function, after this step the sizes of the different parts of the IP prefix are calculated, as well as the header of this memory block in the implemented structure table where the L1 entries will be arrayed. Finally, it sets the table context, it initializes the cache and the tree.

```
1   odph_table_t odph_iplookupv6_table_create(const char *name,
2             uint32_t p1 ODP_UNUSED,
3             uint32_t p2 ODP_UNUSED,
4             uint32_t value_size){
5      odph_iplookup_table_impl *tbl;
6      odp_shm_t shm_tbl;
7      odp_queue_t queue;
8      odp_queue_param_t qparam;
9         .
10        .
11        .
12     tbl = (odph_iplookup_table_impl *)odp_shm_addr(shm_tbl);
13     memset(tbl, 0, impl_size + l1_size);
14
15     tbl->l1e = (prefix_entry_t *)(void *)((char *)tbl + impl_size);
16     for (i = 0; i < ENTRY_NUM_L1; i++)
17        tbl->l1e[i].nexthop = ODP_BUFFER_INVALID;
18     snprintf(tbl->name, sizeof(tbl->name), "%s", name);
19     tbl->magicword = ODPH_IP_LOOKUP_TABLE_MAGIC_WORD;
20     tbl->nexthop_len = value_size;
21
22     for (i = 0; i < 2; i++) {
23        tbl->cache_count[i] = 0;
24        odp_queue_param_init(&qparam);
25        qparam.type = ODP_QUEUE_TYPE_PLAIN;
26        sprintf(queue_name, "%s_%d", name, i);
27        queue = odp_queue_create(queue_name, &qparam);
28        if (queue == ODP_QUEUE_INVALID) {
29           ODPH_DBG("failed to create queue");
30           cache_destroy(tbl);
31           return NULL;}
32        tbl->free_slots[i] = queue;
33        cache_alloc_new_pool(tbl, i);}
34     if (trie_init(tbl) < 0) {
35        odp_shm_free(shm_tbl);
```

```
36        return NULL;}
37      return (odph_table_t)tbl;
38  }
```

<div align="center">Listing 3.6 – Table create function</div>

- **odph_iplookupv6_table_lookup** function performs the table lookup to find any match with the received name.

```
1   odph_iplookupv6_table_lookup(const char *name){
2       odph_iplookup_table_impl *tbl = NULL;
3       odp_shm_t shm;
4       if (name == NULL || strlen(name) >= ODPH_TABLE_NAME_LEN)
5           return NULL;
6       shm = odp_shm_lookup(name);
7       if (shm != ODP_SHM_INVALID)
8           tbl = (odph_iplookup_table_impl *)odp_shm_addr(shm);
9       if (
10          tbl != NULL &&
11          tbl->magicword == ODPH_IP_LOOKUP_TABLE_MAGIC_WORD &&
12          strcmp(tbl->name, name) == 0)
13          return (odph_table_t)tbl;
14      return NULL;
15  }
```

<div align="center">Listing 3.7 – Table lookup function</div>

- **odph_iplookupv6_table_destroy** functions to destroy the table. In this function the correct magic word for the IP prefix table will be checked, then they will call the *trie_destroy* function, after this process will be free all the L2 and L3 entries destroying all L3 subtrees of each L2 subtree and then destroying the L2 subtree, finally the cache is destroyed and the memory is cleared.

```
1   odph_iplookupv6_table_destroy(odph_table_t tbl)
2   {
3       int i, j;
4       odph_iplookup_table_impl *impl = NULL;
5       prefix_entry_t *subtree = NULL;
6       odp_buffer_t *buff1 = NULL, *buff2 = NULL;
7       .
8       .
9       .
10      trie_destroy(impl, impl->trie);
11      buff1 = ENTRY_BUFF_ARR(impl->l1e);
12      for (i = 0; i < ENTRY_NUM_L1; i++) {
13          if ((impl->l1e[i]).child == 0)
14              continue;
15          subtree = (prefix_entry_t *)impl->l1e[i].ptr;
```

```
16        buff2 = ENTRY_BUFF_ARR( subtree );
17        for (j = 0; j < ENTRY_NUM_SUBTREE; j++) {
18          if (subtree[j].child == 0)
19            continue;
20          odp_queue_enq(
21              impl->free_slots[CACHE_TYPE_TRIE],
22              odp_buffer_to_event(buff2[j]));}
23        odp_queue_enq(
24            impl->free_slots[CACHE_TYPE_TRIE],
25            odp_buffer_to_event(buff1[i]));}
26      cache_destroy(impl);
27      odp_shm_free(odp_shm_lookup(impl->name));
28      return 0;
29    }
```

Listing 3.8 – Table destroy function

- ***odph_iplookupv6_table_put_value*** functions to add new values into the table. This function is going to verify if the IP, Key, Prefix, and value are not null. The IP will be parsed once its value is obtained. After this process, are set the L1 entries, the values are inserted into the tree depending on the value of the prefix. If it is less than the root prefix (16 bits) it will be inserted using the function *prefix_insert_into_lx*; this function will handle the last 8 bits when it is inserting into the table. It is used the function *prefix_insert_iter* if the prefix is higher than the root and it will be used to insert groups of octets.

```
1  int odph_iplookupv6_table_put_value(odph_table_t tbl, void *key, void
       *value){
2    odph_iplookup_table_impl *impl = (void *)tbl;
3    odph_iplookupv6_prefix_t *prefix = (odph_iplookupv6_prefix_t *)key;
4    prefix_entry_t *l1e = NULL;
5    odp_buffer_t nexthop;
6    int ret = 0;
7    _uint128_t lkp_ip = 0;
8    if ((tbl == NULL) || (key == NULL) || (value == NULL))
9      return -1;
10   nexthop = *((odp_buffer_t *)value);
11   if (prefix->cidr == 0)
12     return -1;
13   ret = odph_ipv6_addr_parse(&lkp_ip, "ffff:ffff:ffff:ffff:ffff:ffff:
       ffff:ffff");
14   if (ret < 0) {
15     printf("Failed to get IPv6 addr from str\n");
16     return -1;}
17   prefix->ip = prefix->ip & (lkp_ip << (IP_LENGTH - prefix->cidr));
18   ret = trie_insert_node(
19         impl, impl->trie,
```

```
20              prefix->ip, prefix->cidr, nexthop);
21      if (ret < 0) {
22        ODPH_DBG("failed to insert into trie\n");
23        return -1;}
24      l1e = &impl->l1e[prefix->ip >> 112];
25      odp_buffer_t *buff = ENTRY_BUFF_ARR(impl->l1e) + (prefix->ip >> 112)
          ;
26      if (prefix->cidr <= 16) {
27        ret = prefix_insert_into_lx(
28            impl, l1e, prefix->cidr, nexthop, 16);
29      } else {
30        ret = prefix_insert_iter(impl, l1e, buff,
31            ((prefix->ip) << 112), prefix->cidr - 16,
32            nexthop, 24, 2);}
33      return ret;
34    }
```

Listing 3.9 – Table Put value function

- ***odph_iplookupv6_table_get_value*** function gets a value stored into the table.
  This function is going to verify if the table, Key, and buffer are not null, then the
  L1 entry will be obtained and will start searching in the tree, the resulting data will
  be copied, and if there is not a match will return only match the default prefix and
  an ODP buffer invalid.

```
1   int odph_iplookupv6_table_get_value(odph_table_t tbl, void *key,
2           void *buffer ODP_UNUSED,
3           uint32_t buffer_size ODP_UNUSED){
4     odph_iplookup_table_impl *impl = (void *)tbl;
5     _uint128_t ip;
6     prefix_entry_t *entry;
7     odp_buffer_t *buff = (odp_buffer_t *)buffer;
8
9     if ((tbl == NULL) || (key == NULL) || (buffer == NULL))
10      return -EINVAL;
11
12    ip = *((_uint128_t *)key);
13    entry = &impl->l1e[ip >> 112];
14    if (entry == NULL) {
15      ODPH_DBG("failed to get L1 entry.\n");
16      return -1;}
17    ip <<= 112;
18    while (entry->child) {
19      entry = (prefix_entry_t *)entry->ptr;
20      entry += ip >> 24;
21      ip <<= 8;}
22    if (entry->nexthop == ODP_BUFFER_INVALID) {
23      printf("only match the default prefix\n");
```

```
24        *buff = ODP_BUFFER_INVALID;
25      } else {
26        *buff = entry−>nexthop;}
27      return 0;
28  }
```

Listing 3.10 – Table Get Value function

- ***odph_iplookupv6_table_remove_value*** functions to remove a entry from the table. This function will start checking if the table and the key (including the prefix) are not null, then it detects if there is a prefix overlap using the *trie_detect_overlap* function, if the prefix is less than the root it is used the function *prefix_delete_lx*. If the return of this the function is equal to 1, the next $2^8$ entries will equal to over_cidr and over_nexthop. In this case, will be not pushed the over_cidr and over_nexthop to the next level. In the other case when the prefix is higher than the root, it is used the function *prefix_delete_iter* destroying the subtrees, after this, it checks if it can recycle the entry. An entry can be recycled due to two reasons: all children of the entry are the same, or all children of the entry have a prefix smaller than the level bottom bound. At the end the function will finish with the *trie_delete_node*.

```
1   int odph_iplookupv6_table_remove_value(odph_table_t tbl, void *key) {
2     odph_iplookup_table_impl *impl = (void *)tbl;
3     odph_iplookupv6_prefix_t *prefix = (odph_iplookupv6_prefix_t *)key;
4     _uint128_t ip;
5     uint8_t cidr;
6
7     if ((tbl == NULL) || (key == NULL))
8       return −EINVAL;
9     ip   = prefix−>ip;
10    cidr = prefix−>cidr;
11    if (cidr == 0)
12      return −EINVAL;
13    prefix_entry_t *entry = &impl−>l1e[ip >> 112];
14    odp_buffer_t *buff = ENTRY_BUFF_ARR(impl−>l1e) + (ip >> 112);
15    uint8_t over_cidr, ret;
16    odp_buffer_t over_nexthop;
17    trie_detect_overlap(
18        impl−>trie, ip, cidr, 16, &over_cidr, &over_nexthop);
19    if (cidr <= 16) {
20      prefix_delete_lx(
21        impl, entry, buff, cidr, over_cidr, over_nexthop, 16);
22    } else {
23      prefix_entry_t *ne = (prefix_entry_t *)entry−>ptr;
24      odp_buffer_t *nbuff = ENTRY_BUFF_ARR(ne);
25      ne += ((_uint128_t)(ip << 112) >> 24);
```

```
26        nbuff += ((_uint128_t)(ip << 112) >> 24);
27        ret = prefix_delete_iter(impl, ne, nbuff, ip, cidr − 16, 24, 2);
28        if (ret && can_recycle(entry, 16)) {
29          /* destroy subtree */
30          cache_init_buffer(
31            *buff, CACHE_TYPE_SUBTREE,
32            sizeof(prefix_entry_t) * ENTRY_NUM_SUBTREE);
33          odp_queue_enq(
34            impl−>free_slots[CACHE_TYPE_SUBTREE],
35            odp_buffer_to_event(*buff));
36          entry−>child = 0;
37          entry−>cidr = over_cidr;
38          entry−>nexthop = over_nexthop;}}
39      return trie_delete_node(impl, impl−>trie, ip, cidr);
40 }
```

Listing 3.11 – Table Remove value function

In addition to the IP prefix lookup table and binary tree, a function is necessary to parse the IP address. The *odph_ipv6_addr_parse* function reads the IP as a string, and it parses into a vector with the 128 bits. The received string segments it into four sections of 32 bits to parse the IP, and then all that section will be joined at the final IP address variable.

```
1  int odph_ipv6_addr_parse(_uint128_t *ip_addr, const char *str){
2    unsigned byte[ODPH_IPV6ADDR_LEN];
3    int i;
4    _uint128_t p_ip1, p_ip2, p_ip3, p_ip4;
5    memset(byte, 0, sizeof(byte));
6    if (sscanf(str, "%02x%02x:%02x%02x:%02x%02x:%02x%02x:%02x%02x:%02x
       %02x:%02x%02x:%02x%02x",
7        &byte[0], &byte[1], &byte[2], &byte[3], &byte[4], &byte[5], &
       byte[6], &byte[7], &byte[8], &byte[9],
8        &byte[10], &byte[11], &byte[12], &byte[13], &byte[14], &byte
       [15]) != ODPH_IPV6ADDR_LEN)
9      return −1;
10   for (i = 0; i < ODPH_IPV6ADDR_LEN; i++)
11     if (byte[i] > 255)
12       return −1;
13
14   p_ip1 = byte[0] << 24 | byte[1] << 16 | byte[2] << 8 | byte[3];
15   p_ip2 = byte[4] << 24 | byte[5] << 16 | byte[6] << 8 | byte[7];
16   p_ip3 = byte[8] << 24 | byte[9] << 16 | byte[10] << 8 | byte[11];
17   p_ip4 = byte[12] << 24 | byte[13] << 16 | byte[14] << 8 | byte[15];
18   *ip_addr = p_ip1 << 96 | p_ip2 << 64 | p_ip3 << 32 | p_ip4;
19
20   return 0;
```

```
21    }
```

Listing 3.12 – ODP IPv6 Parse

### 3.2.3   MACSAD lookup

For the LPM MACSAD process, the necessary functions are in charge of creating the tables, add values and get values. MACSAD parses the packets and receives the information of IP to perform the LPM match. The controller reads the entries from the trace file and adds into the tables with the prefix information.

- ***table_create*** function will create the match table; it can be an exact table or an LPM table. It is selected depending on the size of the key if it is IPv4 or IPv6. After this process a table LPM lookup will be performed, if the result is not null the table is destroyed (*odph_iplookup_table_destroy*) then the table is created using the ODP function *odph_iplookup_table_create*.

```
1    case LOOKUP_LPM:
2        snprintf(name, sizeof(name), "%s_lpm_%d_%d", t->name, socketid,
     replica_id);
3        if(t->key_size <= 5){ //IPV4
4            if ((tbl = odph_iplookup_table_lookup(name)) != NULL){
5                odph_iplookup_table_destroy(tbl);}
6            // name, capacity, key_size, value size
7            tbl = odph_iplookup_table_create(name, 2, t->key_size, t->
     val_size);
8            if(tbl == NULL) {
9                debug("  ::Table %s creation fail\n", name);
10               exit(0);}}
11           create_ext_table(t, tbl, socketid);
12       else if(t->key_size <= 17){ //IPV6
13           if ((tbl = odph_iplookupv6_table_lookup(name)) != NULL){
14               odph_iplookupv6_table_destroy(tbl);}
15           tbl = odph_iplookupv6_table_create(name, 2, t->key_size, t->
     val_size);
16           if(tbl == NULL) {
17               debug("  ::Table %s ipv6 creation fail\n", name);
18               exit(0);}}
19           create_ext_table(t, tbl, socketid);
```

Listing 3.13 – MACSAD Table Create

- ***lpm_add*** function will add entries inside the LPM table. As a first step, it is used the Key to verify if it is IPv4 or IPv6. In both cases, the logic is the same, just the ODP functions and the size of the variables will be change to the corresponding

case. In this function, the IP is validated and then it is parsed to be added to the LPM table with *odph_iplookupv6_table_put_value*.

```
1   else if(t->key_size <= 17){
2       key[16] = depth;
3       unsigned byte[ODPH_IPV6ADDR_LEN+1];
4       odph_iplookupv6_prefix_t prefix2;
5       for (int i = 0; i < ODPH_IPV6ADDR_LEN; i++)
6           if (key[i] > 255)
7               return;
8       .
9       .
10      .
11      prefix2.ip = p_ip1 << 96 | p_ip2 << 64 |  p_ip3 << 32 | p_ip4;
12      prefix2.cidr = 64;
13      ext->content[ext->size] = copy_to_socket(value, t->val_size+sizeof
        (int), t->socketid);
14      value3 = malloc(t->val_size);
15      memcpy(value3, value, t->val_size);
16      ret = odph_iplookupv6_table_put_value(ext->odp_table, &prefix2, &
        value3);
17      ext->size++;
18      if (ret == -1) {
19          exit(EXIT_FAILURE);}}
```

Listing 3.14 – MACSAD Table Add

- ***lpm_lookup*** function will be in charge to perform the table get. This function matches the Key to detect if it is IPv4 or IPv6. The IP is parsed, and the lookup process is performed using the *odph_iplookupv6_table_get_value*.

```
1   else if(t->key_size <= 17){
2       unsigned byte[ODPH_IPV6ADDR_LEN+1];
3       memset(byte, 0, sizeof(byte));
4       .
5       .
6       .
7       lkp_ip2 = p_ip1 << 96 | p_ip2 << 64 |  p_ip3 << 32 | p_ip4;
8       ret = odph_iplookupv6_table_get_value(ext->odp_table, &lkp_ip2, &
        result, 0);
9       if (ret < 0) {
10          return t->default_val;
```

Listing 3.15 – MACSAD Table lookup

## 3.3 Controller

The controller in an SDN architecture is in charge of controlling the dataplane elements. In the case of MACSAD the controller is going to handle the actions and the information to be filled in the tables defined at the P4 code. The MACSAD Controller code implemented can be found in Appendix C. For the use case of IPv4 and IPv6 the tables defined at the P4 code are: *ipv6_fib_lpm* and *sendout*.

The controller can read and analyze from an input file (Table Trace) to fill the information in the tables. The Table Trace files have a specific format and information depending on the MACSAD use case. In the case of the IPv4 and IPv6 use cases, the information required is destination IP, MAC, and port addresses. The controller will parse the Table Trace file and extract the corresponding information, creating the variables that are going to be used to fill the values in the tables.

```
1  if (11 == sscanf(line , "%02x%02x:%02x%02x:%02x%02x:%02x%02x:%02x%02x:%02x
       %02x:%02x%02x:%02x%02x %x:%x:%x:%x:%x:%x %d",
2          &values_ip[0], &values_ip[1], &values_ip[2], &values_ip[3],
3          &values_ip[4], &values_ip[5], &values_ip[6], &values_ip[7],
4          &values_ip[8], &values_ip[9], &values_ip[10], &values_ip[11],
5          &values_ip[12], &values_ip[13], &values_ip[14], &values_ip[15],
6          &values[0], &values[1], &values[2],
7          &values[3], &values[4], &values[5], &port) ){
8      if (mac_count==MAX_MACS-1){
9        break;}
10     ++mac_count;
11     for( i = 0; i < 6; ++i )
12       macs[mac_count][i] = (uint8_t) values[i];
13     for( i = 0; i < 16; ++i )
14       ips[mac_count][i] = (uint8_t) values_ip[i];
15     portmap[mac_count] = (uint8_t) port;
16   } else {
17     fclose(f);return -1;
18   }
```

Listing 3.16 – Controller Parse trace file

To fill the table with information in the P4 code, the controller uses the functions *fill_ipv6_fib_lpm_table* and *fill_sendout_table* to fill the information from the Trace file into the tables. The *fill_ipv6_fib_lpm_table* function defines the default actions of the P4 table *ipv6_fib_lpm* to *fib_hit_nexthop*, being necessary the parameters of IP, Port, and destination MAC address. Additionally, this function sends the IP to the ODP LPM table. The *fill_sendout_table* function sets the default action to the P4 table *sendout* to *rewrite_src_mac* being necessary to pass the Port and the defined MAC address. Additionally, a sleep of 1000*us* between every addition it is necessary to let the ODP

functions add the entries without missing any information.

```
1   void init() {
2       int i;
3       uint8_t smac[6] = {0xd0, 0x69, 0x0f, 0xa8, 0x39, 0x90};
4       for (i=0;i<=mac_count;++i){
5           fill_ipv6_fib_lpm_table(ips[i], portmap[i], macs[i]);
6           fill_sendout_table(portmap[i], smac);
7                   usleep(1000);}
8   }
```

Listing 3.17 – Controller initialization function

## 3.4   P4 Code

In a P4 code, the headers of the packets are defined at the beginning of the code, in the case of IPv4 and IPv6 the headers are Ethernet and IP. The complete P4 code implemented can be found in Appendix D.

```
1   header_type ethernet_t {
2       fields {
3           dstAddr : 48;
4           srcAddr : 48;
5           etherType : 16;}}
6
7   header_type ipv6_t {
8       fields {
9           version : 4;
10          trafficClass : 8;
11          flowLabel : 20;
12          payloadLen : 16;
13          nextHdr : 8;
14          hopLimit : 8;
15          srcAddr : 128;
16          dstAddr : 128;}}
```

Listing 3.18 – IPv6 P4 Headers

The necessary tables to forward the packets are the *ipv6_fib_lpm* and *sendout.*

- **_ipv6_fib_lpm_** is going to be in charge to call the LPM function when a valid IPv6 packet arrives. This function reads the destination address of the packet and matches with LPM. The actions if it is found a match are the *fib_hit_nexthop* that modifies the Ethernet designation address with the nexthop information, updates the egress port to the corresponding one, and reduces the hopLimit in one, the other action of the function it is *on_miss* that discards the packet.

- **sendout** table is going to read the egress port and perform an exact match, if a match is found, the action will be *rewrite_src_mac* updating the source MAC address of the forwarded packet.

```
1   action on_miss() {}
2
3   action fib_hit_nexthop(dmac, port) {
4     modify_field(ethernet.dstAddr, dmac);
5     modify_field(standard_metadata.egress_port, port);
6     add_to_field(ipv6.hopLimit, -1);}
7
8   table ipv6_fib_lpm {
9     reads {
10      ipv6.dstAddr : lpm;}
11    actions {
12      fib_hit_nexthop;
13      on_miss;}
14    size : 512;}
15
16  action rewrite_src_mac(smac) {
17    modify_field(ethernet.srcAddr, smac);}
18
19  table sendout {
20    reads {
21      standard_metadata.egress_port : exact;}
22    actions {
23      on_miss;
24      rewrite_src_mac;}
25    size : 512;}
26
27  control ingress {
28    apply(ipv6_fib_lpm);
29    apply(sendout);}
```

Listing 3.19 – IPv6 P4 Tables and Actions

The control ingress is the place where the flow is going to be defined, fist it is going to apply the *ipv6_fib_lpm* table and then the *sendout* table. In the use case of MACSAD all the tables are defined at the ingress.

```
1   control ingress {
2     apply(ipv6_fib_lpm);
3     apply(sendout);}
4
5   control egress {}
```

Listing 3.20 – IPv6 P4 Control Ingress and Egress

In Table 4, we compare the use case complexity of the L2 and L3 forwarding with IPv4 and IPv6. The *Parsing* refers to the headers and fields that are parsed at the P4 code, in the case of L3 forwarding the headers are two, Ethernet and IP and the number of fields depend on the protocol (IPv4 or IPv6). *Processing* contains the information of the tables defined, the IPv4 and IPv6 use cases have two, the *ipv6_fib_lpm* and *sendout.* In *Packet Modification* are mentioned the headers that are added or removed from the final packet. *Metadata* is the local information that passes through the tables, in MACSAD is used to pass the egress port of the switch. In *Action Complexity* are summarized the fields and expressions that are modified, in the L3 cases, the fields are Ethernet source and destination address, the Time to Live (TTL) (hop limit for IPv6), and the metadata. The *Lookups* can be Hash (exact) or LPM, the exact match is used to find the egress port and the LPM for the IP lookup.

Comparing the use cases, the IPv6 has 1 LPM table with the key size of 128 bits being higher than the IPv4 case. The Packets fields parsed in the IPv6 are also more than the parsed in the IPv4. On the other hand, the L2 forwarding case has one header with three fields and two hash lookups of 48 bits.

The flow diagram of the P4 code (Generated using P4 Graphs[3]) is speared in the Parse and the Table flows. The Parse flow shows the logic while parsing a packet. In Figure 14b, it is presented the IPv6 Parse, when a packet arrives, the first step is to parse the Ethernet header if the eherType is IPv6 then the IP header is parsed. After this process the tables are executed. The Table flow shows the pipeline of the ingress and egress, in Figure 15b can be seen the order of the applied tables. First, the *ipv6_fib_lpm* and finally the *sendout* table. After this step, the packet is sent to the egress to be forward.

## 3.5   LPM Flow Diagrams

In this subsection, we describe the flow diagram of the LPM implementation containing the relationship between the implemented function of the binary tree, ODP lookup, and MACSAD lookup. Additionally, it is described the flow diagram of the P4 implementation.

In Figure 16, the main OPD LPM lookup flow is observed, when an ODP program starts, the first steps are to initialize the ODP global and local variables (Internal functions of ODP), with this function are set parameters as the distribution of cores, the RX and TX ques, the memory blocs to be used. The tables (Figure 17) are created after this step. Once the tables are created, if new action is requested the flow will continue, the table actions (Figure 18) as described in the previous subsections, the Table Put values, Tables Get Values, Table Remove Value, and the Tables Destroy (Figure 17).

---

[3]   <https://github.com/p4lang/p4-hlir/blob/master/bin/p4-graphs>

(a) IPv4 Parse          (b) IPv6 Parse

Figure 14 – L3-FWD Parse Flow



(a) IPv4 Table          (b) IPv6 Table

Figure 15 – L3-FWD Table Flow

Figure 16 – ODP LPM main flow

The MACSAD LPM implementation flow diagram is described in Figure 19. Following the same process as described above when an ODP program starts, the OPD is initialized, then the tables are created (Figure 20). MACSAD has two types of tables, the Lookup Exact, and the Lookup LPM. The exact match for the L2 information is used (Source and destination MAC address), and available types are the Cuckoo and the Hash

table.

MACSAD uses a key to pass the IP and prefix values to the functions. For IPv4 the key size is 5, where the first 4 are the IP, and the last one is the prefix that will be used for the LPM table. For the IPv6 case, the size of the key is 17 where the first 16 are the IP and the last one the prefix. For the LPM lookup tables, the code will match the size of the key; if it is 5, it will be an IPv4 address creating the corresponding table, if the key is 17, the table created will be an IPv6 type. For the MACSAD use case, one table is created for the LPM values. The implemented table actions (Figure 21) are the Table Add and Table Get. The Table Add matches the size of the key with the same logic of the table create and the corresponding IPv4 and IPv6 ODP put value functions call (*odph_iplookup_table_put_value* or *odph_iplookupv6_table_put_value*). The Table Get uses the same logic for the size of the key and depending on the size will be called the the ODP get value functions (*odph_iplookup_table_get_value* or *odph_iplookupv6_table_get_value*).

## 3.6   Summary

In this chapter we described the IPv4 LPM support implementation in ODP including our MACSAD support for the LPM library, we detailed the IPv6 implementation including our LPM helper library in ODP and our MACSAD support that coverage the controller and the P4 Code design. In the ODP implementation is described the different functions that are part of the LPM process, presenting our IPv6 LPM implementation and comparing it with the IPv4 library developed by ODP. We illustrated the complete workflow of the LPM implementation including the MACSAD and ODP library parts, having a comprehensive view of the process and how the different projects are integrated, also comparing the use case complexity of the L2 and L3 forwarding with IPv4 and IPv6.

In the next chapter, we are going to present our Packet Generator Tool BB-Gen, that will support us in the evaluation of the implemented functions and libraries, allowing us to generate the different packet traces for the performances test.

Table 4 – P4 Use Case Complexity

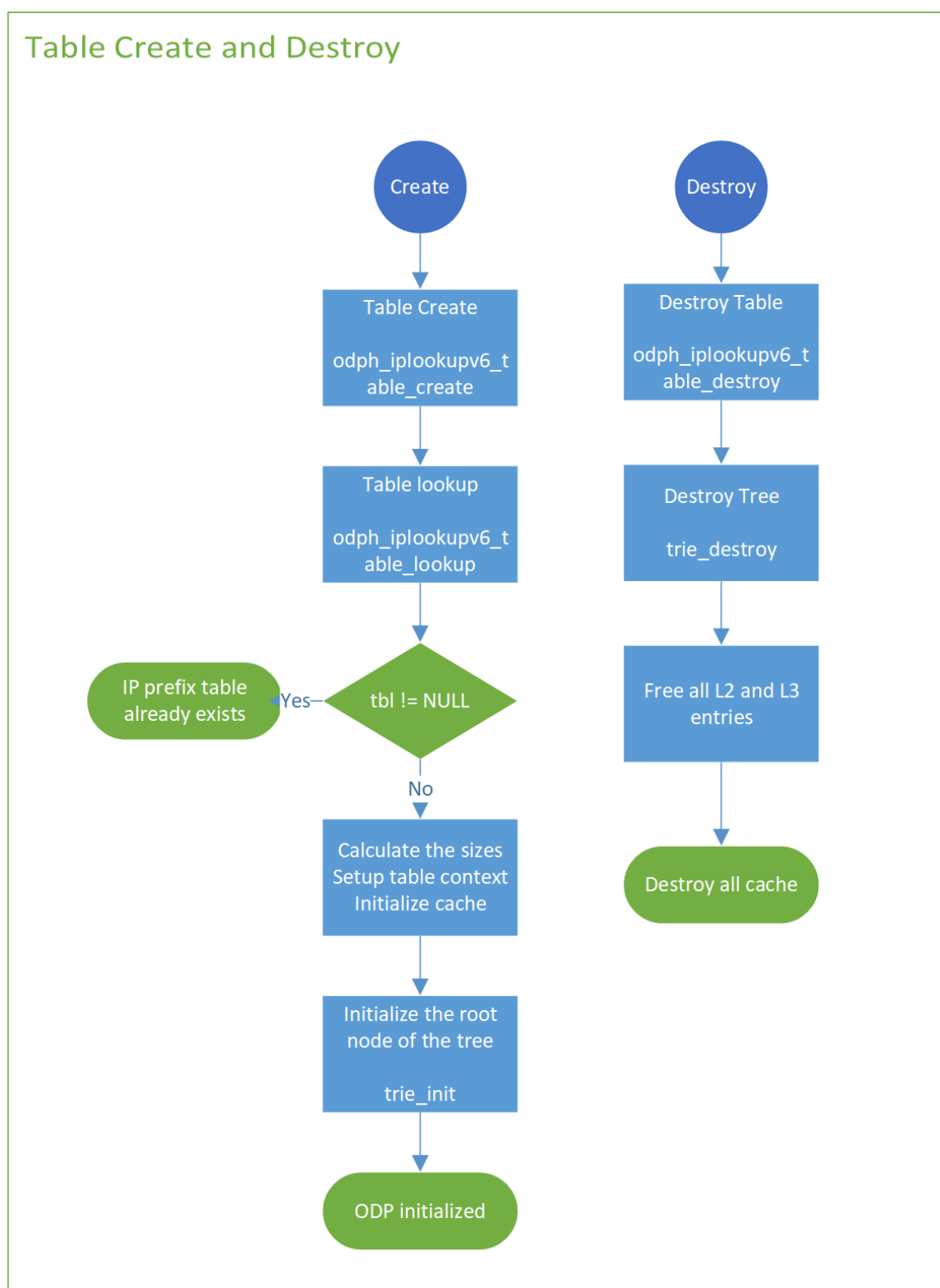| | | L2-FWD | L3-FWD (IPv4) | L3-FWD (IPv6) |
|---|---|---|---|---|
| | | P4_14 | P4_14 | P4_14 |
| Parsing | #Packet headers | 1 | 2 | 2 |
| | #Packet fields | 3 | 13 | 19 |
| | #Branches in parse graph | 1 | 2 | 2 |
| Processing | #Tables (no dep) | 2 | 2 | 2 |
| | Depth of pipeline | 2 | 2 | 2 |
| | Checksum on/off | off | off | off |
| | Table size | 8192 | 512 | 512 |
| State Accesses | #Writes to different register | 0 | 0 | 0 |
| | #Writes to same register | 0 | 0 | 0 |
| | #Reads to different register | 0 | 0 | 0 |
| | #Reads to same register | 0 | 0 | 0 |
| Packet Modification | #Header adds | 0 | 0 | 0 |
| | #Header removes | 0 | 0 | 0 |
| Metadata | #Metadatas | 1 | 1 | 1 |
| | Metadata size(bits) | 9 | 9 | 9 |
| Action Complexity | #Field writes | 2 | 4 | 4 |
| | #Arithmetic expressions | 0 | 0 | 0 |
| | #Boolean expressions | 0 | 0 | 0 |
| Lookups | #Hash_lookups [key_length(bits)] | 2 [48] | 1 [9] | 1 [9] |
| | #LPM [key_length(bits)] | 0 | 1 [32] | 1 [128] |

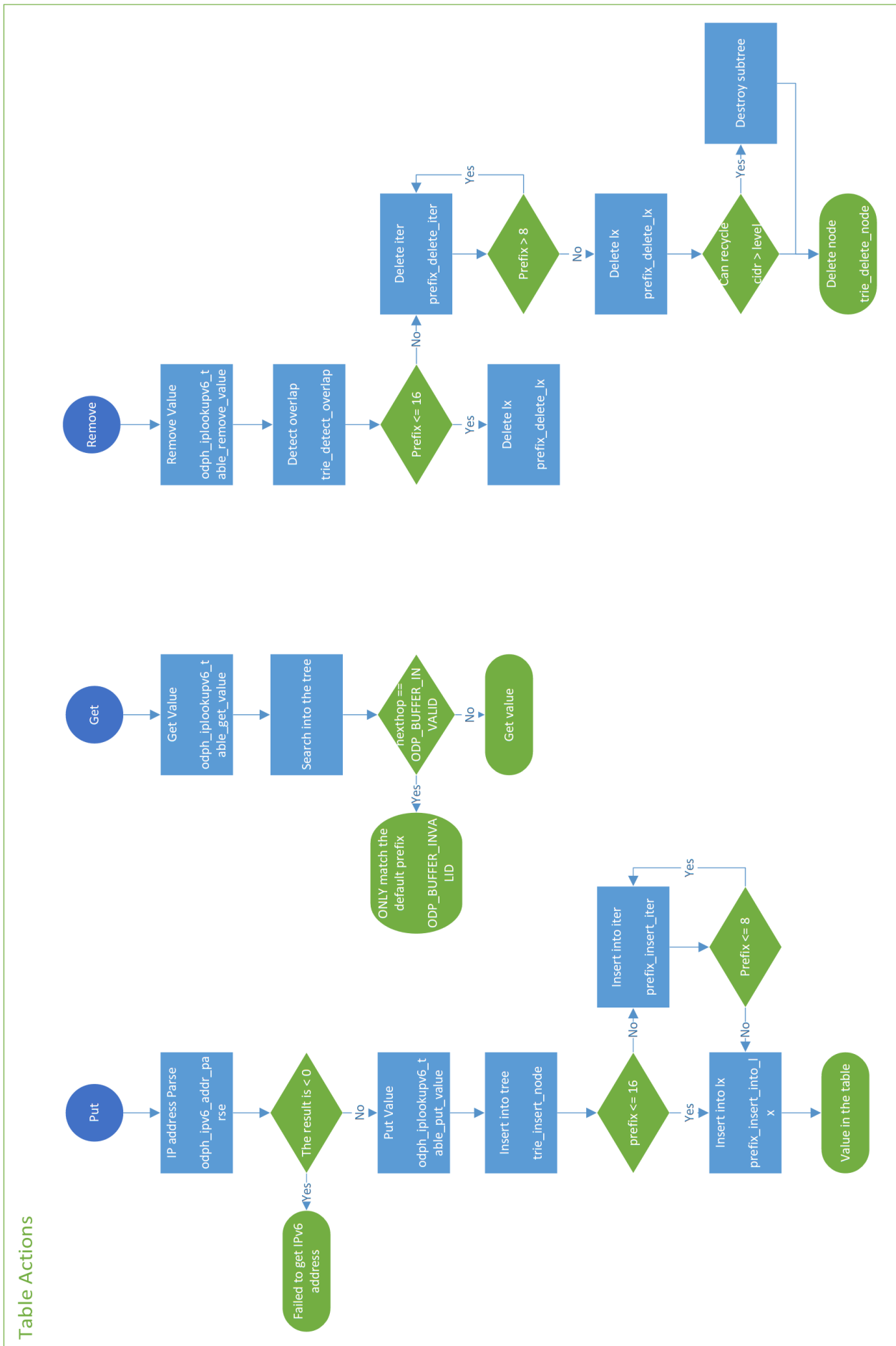Figure 17 – ODP Table Create and Destroy
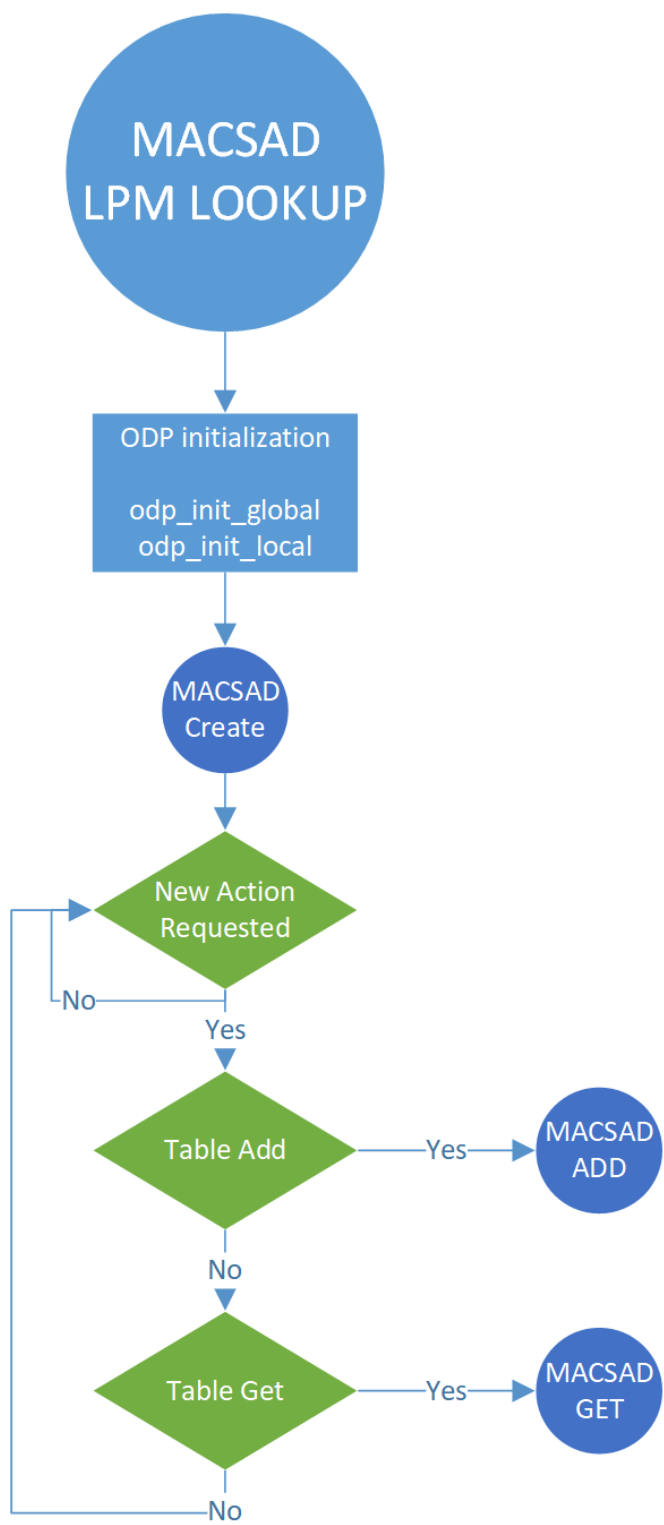
Figure 18 – ODP Table Actions
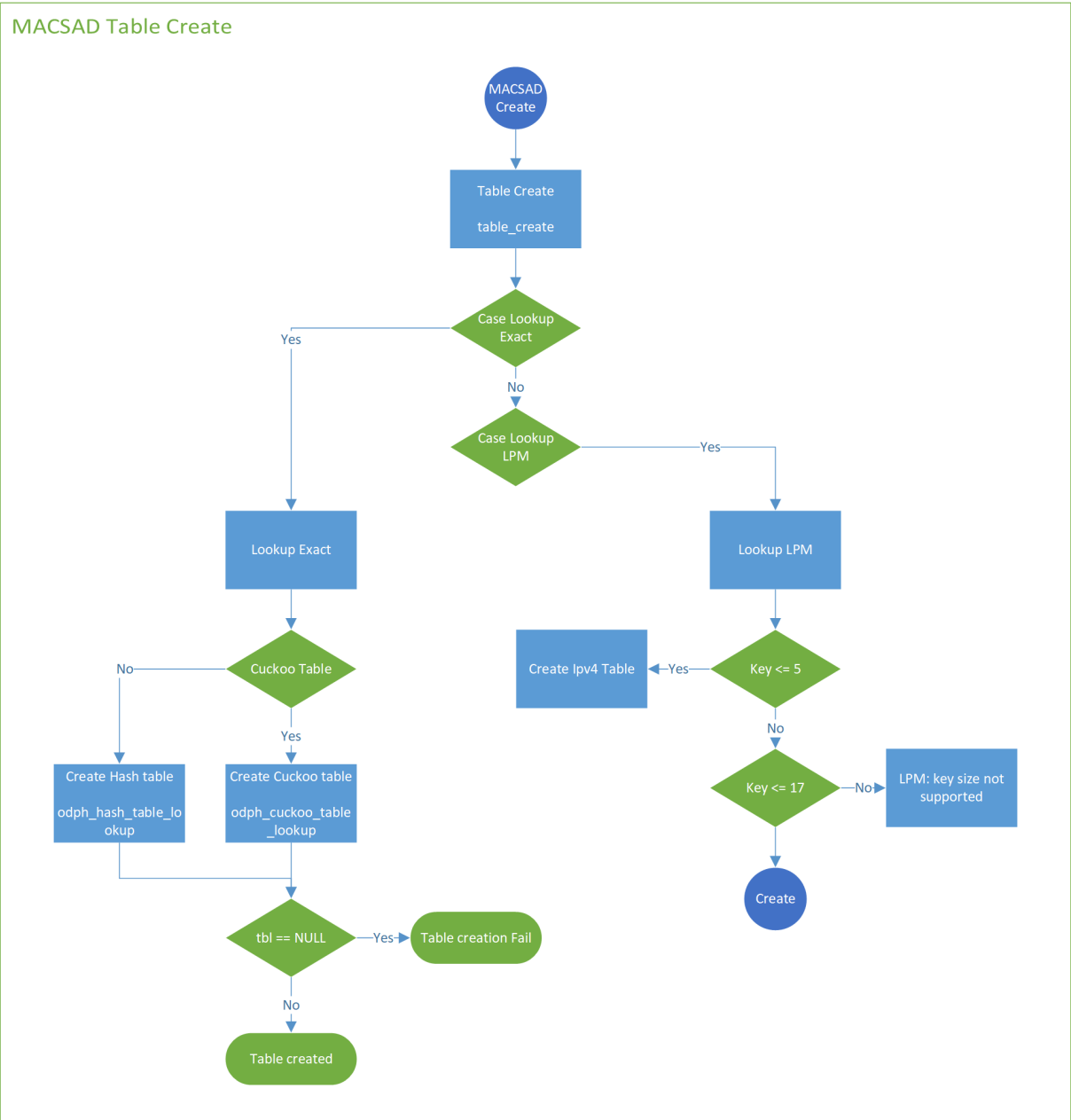
Figure 19 – ODP LPM main flow

Figure 20 – ODP Table Create

Figure 21 – ODP Table Actios

# 4 Packet Generator Tool: BB-Gen

With P4 gaining traction to define datapath pipelines along auto-generated control plane APIs, the protocol-independence and increased flexibility add non-trivial hazards when it comes to functional and in-depth performance evaluation. P4-dependent workload traces are needed along automated methods to populate the tables of the datapath under test accordingly. Without proper tools, manual efforts are required for tedious tasks such as creating appropriate PCAP traces, defining the distribution of field values, and inserting entries in the pipeline tables. To this end, we developed BB-Gen (RODRIGUEZ *et al.*, 2018; CESEN *et al.*, 2018), a packet crafter and table generator tool that given a P4 application and a corresponding user configuration results in packet and table traces to carry automated performance evaluation tasks. We evaluated BB-Gen with P4 applications of increasing complexity (from L2 to VXLAN-based Data Center Gateway), using two different multi-architecture backend compilers (MACSAD, T4P4S) and different targets.

With increasing number of services over the Internet like email, web, video streaming and so forth, the demand for bandwidth is increasing exponentially. Along with it, the necessity to evaluate and test network capabilities become prevalent. While networks are becoming more (re-configurable), network testing tools are becoming equally complex adapting to the need of the hour. The network testing and benchmarking tools depend on network workload generation to simulate the network traffic for testing purposes. This trivial task has been the foundation for several research activities like (BOTTA *et al.*, 2012) focusing towards performance, scalability, and reliability of networks and network devices.

Traffic generator tools are an essential part of network testing with features ranging from supporting list of protocols, analyzing network traffic or measuring throughput to calculating latency of packets. In theirs strive to achieve feature completeness, the tools are getting more complex each time, and making it hard to port, manage, and use. To address this, we propose BB-Gen which is a python based tool, with a primary focus on simplicity, excelling in the creation of network packet traces.

BB-Gen is a simple CLI based packet crafter written in Python over Scapy library. It can natively craft packets for different standard and custom protocols. It aims to create PCAP files to be used with a wide set of Traffic Generators (e.g., pktgen-dpdk (OLSSON, 2005), NFPA (CSIKOR *et al.*, 2015a; CSIKOR *et al.*, 2015b), TCPDUMP (TCPDUMP, 2010)) helping network developers to validate the network and execute performance tests over the targets.

Though BB-gen is primarily used to create PCAP trace files, it differs itself by generating the table trace files for the PCAPs which are necessary to fill the table flow configuration of the target device for the network testing. Table traces contains the main information of the generated packets (e.g. source/destination IP/MAC address). BB-Gen allows to create traces files with same/random IP/MAC/L4Port details showing its control over the header fields like source and destination MAC addresses, IP addresses, TCP or UDP ports while creating packets. It also allows the user to create a complete set of PCAPs for performance test by specifying a single flag in the command line. Under this performance setting, PCAPs generated comprises of all the standard packet sizes (64, 128, 256, 512, 1024, 1280, 1518) (BRADNER; MCQUAID, 1999) and also features simple (best-case) and complex (worst-case) scenarios by using typical/random distribution sets of header fields respectively. A single command can generate both PCAP and table trace files. The command line arguments which are reckoned and self explained go in agreement with the easy use of BB-Gen. A custom protocol support to BB-Gen can be easily added by first adding the support to Scapy similar to the Contrib[1] and then extending BB-Gen protocol list with minimal code changes.

## 4.1   Architecture

Figure 22 shows the principal components of the architecture of BB-Gen Packet Crafter.



Figure 22 – BB-Gen Architecture

---

[1]   https://github.com/secdev/scapy/tree/master/scapy/contrib

- **User:** The user introduces required parameters such as distribution, protocols, numbers of entries, use case[2], packet sizes, necessary to create the trace files. [3]

- **Core:** Being the principal part of BB-Gen, it receives and process information from the *User*, and generates the packet details to be included in the trace's files. It comprises of three sub-modules i.e., Parser, Data Generator, and Packet Generator as explained below (Fig. 23).

Figure 23 – BB-Gen Core module

– **Parser:** it is in charge of selecting the protocols to be used as well as the distribution, using the information introduced by the use or the default values in case of missing information (Fig. 24).

Figure 24 – BB-Gen Parser module

– **Data Generator:** using the protocols and the distribution details from *Parser*, it generates the list of source and destination IP, MAC and Ports (Fig. 25).

– **Packet Creator:** with the information set at *Parser* and the list of IP, MAC and Ports generated at *Data Generator*, the Packet Creator is going to create the list of packets with all the defined fields. With the list of packets prepared, the table trace file is going to be created using the informations about packet

---

[2]   Supported use case: MACSAD
[3]   More information on required parameters to generate the traces are described in BB-Gen GitHub Wiki page, <https://github.com/intrig-unicamp/BB-Gen>

Figure 25 – BB-Gen Data Generator module

contents. And finally, the list of packets is sent to the Scapy block to generate the final PCAP trace (Fig. 26).



Figure 26 – BB-Gen Packet Creator module

- **Scapy:** is composed of the *Packet* and *PCAP* sub-blocks. The *Packet* is going to assemble the packets included in the list of packets with the correct protocol format. The *PCAP* will read the assembled packets and generate the PCAP file completing the BB-Gen process (Fig. 27).



Figure 27 – BB-Gen Scapy module

- **Output Files:** BB-Gen generates two output files, the Table trace and the PCAP trace:

  - **Table Trace:** Generated by the Packet Creator module, it is a plain text containing matching entries derived from the PCAP packet trace and Packet List. It can be customized to the specific use-case and used to populate entries to the dataplane tables.

  - **PCAP Trace:** Generated by the Scapy module with all the information and protocols configured by the user. It can be used with a wide set of benchmarking tools for the evaluation test.

## 4.2   Main features

The principal features and capabilities of BB-Gen are summarized as below:

- Designed for simplicity, BB-Gen delivers an intuitive CLI based interface. By specifying only a few flags, the user can create a set of traces files.

- A P4 code can be used as an Input. BB-Gen identifies the protocols and auto-generate the traffic traces.

- Very useful for best-case and worst-case testing. It allows to specify a simple/random distribution of header fields sufficient to address the most complex test cases.

- Being a python based tool, it is easy to build, use and extend to support additional protocols and new features.

- Easily create multiple PCAPs in a single step. The user can define the number of flows, packet sizes, for each PCAP.

- Generates table trace files along with every set of PCAPs utilizing the informations from the PCAP files such as list of IP addresses, MAC address, Port numbers and also the packet encapsulation data for protocols like VXLAN and GRE. Trace file generation is seamless and does not require any additional user input.

- For scalability testing purposes, it can generate traces with more than 1 million unique packet details.

- Supports a list of common standard protocols:

  - Ethernet.
  - IPv4, IPv6.
  - TCP, UDP.
  - Protocol Encapsulations such as GRE and VXLAN.

- Useful for performance tests as it can automatically create packets of different sizes according to the RFC 2544 (BRADNER; MCQUAID, 1999) (64, 128, 256, 512, 1024, 1280, 1518 Bytes) by setting a single performance flag in CLI.

- User defined custom packet sizes are also accepted at the CLI, just being limited by the defined minimum protocol size.

- Accepts user defined payload information. For this scenario, minimum packet size is maintained to be 64 Bytes by padding with random strings if necessary. In case payload saturates the 64 Bytes, the packet size is determined by the aforementioned payload.

Figure 28 – BB-Gen integration with NFPA and MACSAD & T4P4S

- The generated PCAP trace files are accepted as *inputs* for different network benchmarking and performance tools.

- It is a cross-platform tool with support for Windows, Linux, BSD and Mac OS X platforms.

- It is an open source project following BSD 3-Clause License.

## 4.3   Use Case

In order to demonstrate the usability of BB-Gen, we present a use case featuring a programmable dataplane (MACSAD & T4P4S) and a network performance evaluation tool (NFPA) which accepts a set of PCAP trace files as input for each 'determined setup'/'specified configuration' experiment run.

### 4.3.1   Evaluation

While using the NFPA benchmarking tool for performance evaluation of MACSAD & T4P4S, BB-Gen generates the necessary PCAP and table trace files for worst-case scenarios with random header field values (MAC and IP addresses, Port numbers.) from the P4 file set as an input. Multiple sets of PCAPs are created for different packet sizes according to the RFC 2544, and also with a different number of packet flows (100 to 1 million unique flows). This is repeated for each use-cases supported by MACSAD & T4P4S such as L2-Fwd, L3-Fwd with IPv4 and IPv6, GRE, and VXLAN.

For the use case evaluation, the NFPA standalone node is connected to the DUT (MACSAD or T4P4S) as presented in Figure 28 (A detailed use case is presented in Chapter 5). The user defines his P4 program and configures the benchmarking tool accordingly

(in our case, NFPA only requires the high-level configuration details, e.g., packet size). Then, the P4 program is fed into BB-gen. Users can define target PCAPs with different packet sizes (i.e., from 64 to 1500 Bytes) from best-case (i.e., fixed, single header fields) to worstcase (random, unique field values) workloads. While the packets are generated, the DUT is compiled by our Target Compiler module supporting multi-architecture P4 compilers, such as MACSAD and T4P4S.

Once the DUT is running, BB-gen uses the available APIs[4] to carefully populate the flow tables according to the P4 Table Trace containing matching entries derived from the PCAP packet trace and Packet List. At the same time, the BB-gen loads the generated PCAP file into NFPA, which takes care of the practical measurement conforming the standards.

## 4.4   Summary

In this Chapter, we presented BB-Gen a packet generator tool that can generate a suitable trace file used for performance evaluation with a simple process. BB-Gen can also create multiple sets of trace files with different packet sizes and different flow distribution simultaneously along with the corresponding table trace file for each PCAPs and pipelines. The integration with an extensive set of benchmarking tools reinforces the ease of use of the tool and benefits for the community with the evaluations.

We showed how BB-Gen, NFPA and "MACSAD & T4P4S" trio can work conjointly to exemplify how a P4 program suffices to define the datapath pipeline, create match + action table traces to control the P4 DUT, and generate a trace file for the evaluation, covering different complexities and requirements of the project.

---

[4]   As today, APIs are target-specific but P4Runtime (PRT) API support is underway.

# 5 Experimental Evaluation

In this chapter, we evaluate the performance of the two LPM use cases using three different packets I/O engines (DPDK, Netmap, and Socket_mmap). For each combination, we explore the scalability for different workloads (packet traces, table entries) and configuration options (e.g., CPU cores) using NFPA as a benchmarking tool. To generate the traces we use our packet crafter tool BB-Gen[1] that will provide the necessary PCAP files to be used with NFPA. The pipeline implementation[2] and other information for reproducibility purposes including the P4 programs[3] used by MACSAD and the traffic generator tool (BB-Gen) are publicly available.

## 5.1 Testbed and Methodology

Our testbed (Fig. 29) contains two Lenovo ThinkServer RD640 servers with Intel Xeon E5-2620v2, 6 Cores, Hyper-Threading disabled, running at 2.1GHz, 8*8GB DDR3, a dual-port Intel X540-AT2 NIC (10G), and run with Ubuntu Linux 16.04 LTS (kernel 4.4). The Tester server (module 5 in Fig. 29) runs NFPA with DPDK v17.08 and PktGen v3.4.5, and it is connected to the DUT (BRADNER; MCQUAID, 1999) (module 4). The DUT supports multiple packets I/Os to illustrate the ability to accommodate different platform features, such as DPDK v17.08, ODP v1.16.0.0, Netmap v11.2, and the basic Linux Socket_mmap provided by the Linux kernel. The P4 Code (module 1) with the Ipv4 or IPv6 implementations is the first input of the test. The code will be used by BB-Gen (module 2) to auto-generate the PCAP and the table Trace, and by the compiler (module 3) to generate the MacS using the ODP APIs. The controller parses the table Traces file and populated the information into the corresponding table. The PCAP files are used by the NFPA to generate the traffic for the performance evaluation. The MacS is configured to forward packets received from one port to the other and eventually back towards NFPA, which in turn analyzes the packet throughput concerning PPS and BPS (module 6).

For both L3-IPv4 and L3-IPv6, different number of cores (1, 2, 4, and 6) are allocated to the DUT, distinctive workloads are configured by setting different number of IP prefixes (100, 1K, 10K, 100K, 1M) in the lookup table and a matching number of L3 flows used in the synthetic traces.

---

[1] <https://github.com/intrig-unicamp/BB-Gen>
[2] <https://github.com/intrig-unicamp/macsad>
[3] <https://github.com/intrig-unicamp/macsad-usecases>

Figure 29 – Testbed

## 5.2 Performance Analysis

### 5.2.1 L3-IPv4

Figure 30 shows the performance of L3-IPv4 for different FIB sizes and packet I/O drivers, i.e., DPDK, Netmap, and Socket_mmap. The red horizontal lines refer the line rate for different packet sizes. (i.e., 8.44 Mpps for 128 bytes and 4.52 Mpps for 256 bytes). The results for L3-IPv4 are grouped into three sectors indicating different packet sizes (i.e., 64, 128, 256). Each sector is further divided into five different points marking the complexity of the pipeline, i.e., the size of the FIB (e.g., 100, 1K). It can be observed that MacS, with DPDK, reaches the 256 bytes' packet size line rate regardless of the FIB table size. The performance of Netmap is comparatively lower, but it reaches line rate with 512B packets. Also, it is interesting to note that, the measured results for 1M FIB entries are better than for 100K FIB entries. From the results, it is clear that the Linux Socket_mmap driver never saturates the 10G interfaces even with the largest packets (1518 bytes) due to the highly increased number of system calls, fundamental kernel scheduling, and costly context switching imposed by the Linux kernel itself.

In a configuration using 4 cores (Figure 31) the results of the different packet I/O drivers where grouped depending on the FIB size and the packet size. It is important to notice that all the results independent of the driver are starting from '0' in the '$y$' axis. It is clear that the performance increases notably compared to the 1 core setup. With DPDK the line rate is reached even with small packets of 64 bytes and with all the FIB sizes. Netpmap also increases the performance and can saturate the link with packets of 256 bytes. It is interesting to note that the performance results with 64 bytes'

Figure 30 – IPv4 forwarding performance for different I/O drivers (1 CPU core)



Figure 31 – IPv4 forwarding performance for different I/O drivers (4 CPU core)

packets and 1K entries are better than the 100 entries results. This behavior can be caused by a sub-optimal synchronization and CPU utilization of Netmap, it uses NIC interrupts and standard kernel synchronization mechanisms to block on empty or full NIC queues (LETTIERI *et al.*, 2017). With Socket_mmap driver are archived slightly superior results if we compare with the performance of the 1 core evaluation using the same driver.

Figure 32 presents the results for different core configurations (i.e., from 1 to 6), and a 64 bytes' packet size. The results include a FIB size of 100 and 10K. Our results demonstrated that the performance increase with the addition of cores, this analysis applies for the DPDK and Netmap drives, with both configurations and independent of the FIB size it is clear that throughput raise. In case of DPDK, the line rate is reached with 4 cores, as it was demonstrated before, with Netmap the link is saturated only with packets starting from 256 bytes. A similar pattern of the previous results was obtained with Socket_mmap, confirming that the limitations of the Linux kernel affect the performance using this driver.

Figure 32 – IPv4 different cores performance (64 bytes packets)



Figure 33 – IPv6 forwarding performance for different I/O drivers (1 CPU core)

## 5.2.2   L3-IPv6

Results for the L3-IPv6 use case with 1 CPU core are shown in Figure 33. The performance results lead to a similar conclusion to L3-IPv4 where DPDK reaches line rate (red line) with 256 bytes' packets for all FIB sizes. There are some performance differences in case of the Netmap driver, a slight drop as the number of FIB entries grows. However, when comparing our results to the L3-IPv4, we must point out that the peculiarity with 100K and 1M number of entries observed before also applies for L3-IPv6.

Comparing the IPv6 results with IPv4 (Figure 34, it is important to highlight the fact that the performance in the case of IPv4 is slightly superior with different packets sizes (i.e., 64 and 128 bytes). Additionally, it is remarkable that in the case of IPv6 the Key size is 128 bits increasing the complexity of the LPM tree.

The IPv6 performance results with a four cores configuration are presented in Figure 35. It is important to highlight, with Netmap when the FIB size reaches 1M the line rate is not achieved even with the largest packet size, the limitations of Netmap can
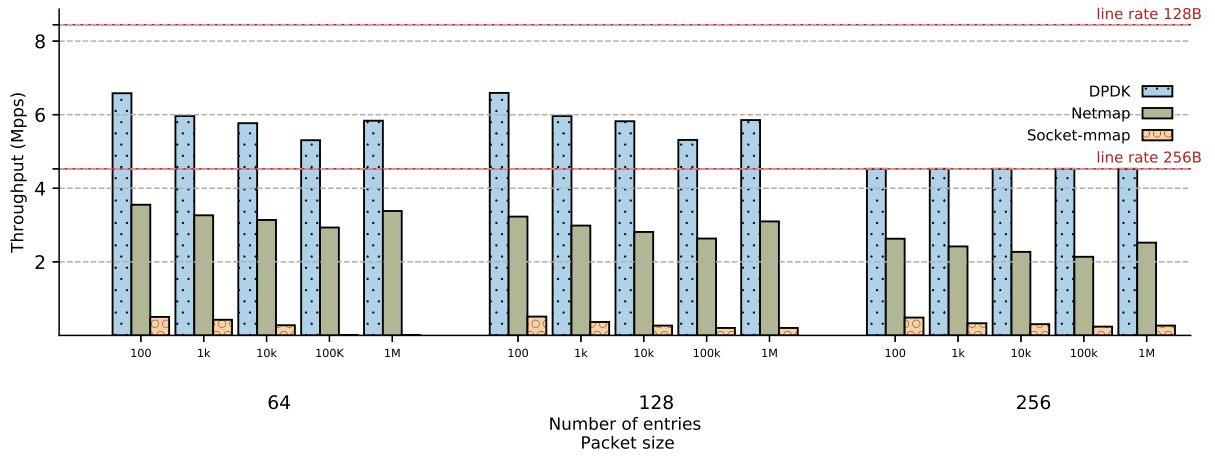
Figure 34 – IPv4/IPv6 forwarding performance for different I/O drivers (1 CPU core).



Figure 35 – IPv6 forwarding performance for different I/O drivers (4 CPU core)

cause this behavior.

Figure 36 shows a throughput comparison as the number of cores increases from 1 to 6. The results with different cores are consistent with what has been discussed before. When the number of CPU cores increases, MACSAD can process more packets resulting in higher throughput.

It is notable (red line in the Figure 37) that as the number of FIB entries increases, the throughput reduces slightly. Moreover, when the table key size increase from 32 bytes (IPv4) to 128 bytes (IPv6) the performance also decreases. This is a significant finding in the understanding of how the complexity of the number of FIB entries and key sizes affect the throughput.

Figure 38 presents the performance results of the IPv6 lookup by setting different values to the prefix length. It is interesting to note that, when the key value (prefix) increases, the performance reduces. With a key of 16 and 32, the line rate is easily reached. When we increase the prefix to 64, the throughput decreases to 12.9 Mpps, and while

Figure 36 – IPv6 different cores performance (64 bytes packets)



Figure 37 – IPv4/IPv6 different cores performance (64 bytes packets)

increases the prefix the performance also reduces. This implies that the throughput is associated with the size of the key. In ODP library when we increase the value of the prefix, the algorithm has to create more leaves, and when we perform a lookup, it is necessary to search in all the levels adding processing and memory consumption.

Comparing MACSAD IPv4 and IPv6 implementation with T4P4s, we can observe in Figure 39 that the performance of the IPv4 use case is slightly superior in MACSAD. This can be caused by an optimized algorithm to perform the LPM lookup. While comparing the IPv4 and IPv6 implementation with the l2fwd use case, we can observe that the L2 forwarding has the worst performance result. From these results, it is clear that exact lookup used by the L2 use case is not optimized what causes a reduction in the performance.

Figure 38 – MACSAD IPv6 prefix length comparison (64 bytes packets, 100 entries)



Figure 39 – MACSAD vs T4P4S use cases comparison (64 bytes packets, 100 entries)

## 5.2.3  CPU Cycles

Performance counters are special hardware registers that are available on most of the CPUs and count the number of some hardware events (i.e., instructions executed, cache-misses suffered, branches mispredicted (Without slowing down the kernel or applications). The perf command (MELO, 2010) is a powerful tool that can instrument CPU performance counters, tracepoints, kprobes, and uprobes (dynamic tracing). Other similar tools to analyze the CPU performance are: the Berkeley Packet Filter (BPF) (MC-CANNE; JACOBSON, 1993), a register-based filter evaluator for filtering network packets (best known for its use in tcpdump), and the Extended BPF (εBPF) variant that analyzes all over the kernel, including maps and used for aggregating statistics of events.

We evaluate the CPU Cycles using Perf and with the same topology of the performance analysis (Fig. 29). For the evaluation, we use a configuration of 1 CPU core and 100 entries, the prefix size of 24, and DPDK as the packet I/O due to the best results during the performance tests. We compare MACSAD and T4P4S; both projects use a similar structure with the difference that T4P4S uses a DPDK backend and libraries,

Figure 40 – MACSAD IPv4 CPU Cycles (324k total cycles, 1 Core, 100 entries)

including the LPM functions.

In Figure 40 it can be seen the results of the IPv4 test, are considered for the graph the results related to the forwarding process. During the test were captured 324k CPU cycles in total. It is important to remark that the commands that process the receiving packets (i.e., *dpdk_recv*, *odp_pktin_recv_tmo*, *ixgbe_recv_pkts_vec*) consume the largest number of cycles (155k CPU cycles). The *table_ipv4_fib_lpm_key* command start the LPM processing using 21k CPU cycles. The MACSAD LPM function of lookup *lpm_lookup* takes 689 CPU cycles samples, inside this command the ODP get function *odph_iplookup_table_get_value* takes 511 CPU cycles. The ODP table create *odph_iplookup_table_create* and table add use few CPU cycles (less than 200). Comparing the LPM commands with the cuckoo commands, it is clear that the LPM functions have a high impact in the CPU cycles consumed.

For the IPv6 case (Figure 41), the results lead to a similar conclusion of the IPv4 results, where the commands that process receiving packets consume more cycles, from the 261k captured more than 100k correspond to the receiving process. In the case of IPv6, the LPM functions consume more CPU cycles; the LPM lookup reaches the 2.6k CPU cycles. The results confirm that the size of the key affects the processing of the packets adding more cycles. It leads to good results, even if the processing in IPv6 is higher, confirming the results from our evaluation of different cores and use cases IPv4/IPv6 where we saw a reduction on the performance when the size of the key increases.

When comparing our results with T4P4S (Figure 42), it worth mentioning that MACSAD uses less CPU cycles. In the case of our LPM lookup, MACSAD IPv4 takes 689 CPU cycles, MACSAD IPv6 2.6k, and T4P4S with DPDK LPM lookup takes 12.9k. With the evaluation, it is possible to see significantly better results for MACSAD using the our implemented ODP library for the LPM process than T4P4S using DPDK library

Figure 41 – MACSAD IPv6 CPU Cycles (261K total cycles, 1 Core, 100 entries)



Figure 42 – MACSAD IPv4/IPv6 vs T4P4S - CPU Cycles - LPM and exact lookup commands (1 Core, 100 entries)

for the LPM process. One difference between our ODP LPM library and DPDK LPM library can be attributed to an optimized algorithm to store the LPM entries in the table, which can reduce the number of cycles consumed by the process. This is consistent with what has been found in our performance evaluation (Figure 39).

## 5.3  Summary

The performance evaluation is a crucial part of the development of a system. We evaluated our LPM implementation with MACSAD in different scenarios and configurations, including different sets of packet I/O drivers, cores, numbers of entries, packet sizes, and much more. We compared our results with other similar projects like T4P4S, analyzing the performance and the CPU Cycles consumed by the projects. We presented

an analysis of how the prefix length can impact the performance and how the size of the prefix tree (number of levels) also impacts on the results.

# 6 Conclusions and Future Work

The focus of this dissertation is the design, implementation, and evaluation of IPv4/IPv6 Longest Prefix Match support in Multi-Architecture Programmable Data-planes to address the limitation of IPv4 lookup support, lack of IPv6 support in ODP, and understanding the performance capabilities. To this end, we have: (i) implemented and evaluated the IPv4 and IPv6 LPM support in MACSAD, (ii) developed a new ODP Helper library for the IPv6 lookup mechanism based on the current IPv4 solution, (iii) carried performance and experimental evaluation of scalability, (iv) developed a P4 CLI-based packet crafter to generate packet flows formatted as PCAP files, and (v) released all artifacts as open source.

This dissertation describes the IPv4 LPM support implementation and a detailed IPv6 implementation including the functions and processes, including the Controller and the P4 Code design. We compared the complexity of three different P4 use cases, i.e., L2-FWD, L3-FWD with IPv4 and IPv6, remarking the different sizes of headers, tables, and lookup types that impact performance evaluation. We learned that the performance is affected by the lookup type and the size of the key, which are key characteristics of the use case pipeline complexity (e.g., number of tables, key size, table size).

We evaluated the performance of the LPM use cases using three different packet I/O engines (DPDK, Netmap, Socket_mmap). We explored the scalability for different workloads (packet traces, table entries) and configuration options (e.g., CPU cores), using NFPA to transmit the PCAPs created by BB-Gen. Comparing the performance of the different packet I/O drivers, we confirmed that Socket-mmap is slower as it is the Linux default driver without any fastpath advantages from Netmap or DPDK. Netmap performance is lower than DPDK on ODP system because in the case of the Netmap packet the copy operation becomes costlier as in the case of the DPDK support, the ODP implement a zero-copy feature. With our findings, we can understand that with the use of the Netmap driver we have a sub-optimal CPU synchronization and CPU utilization, the reason is that Netmap uses NIC interrupts and standard kernel synchronization mechanisms to block on empty or full NIC queues. In conclusion, DPDK seems to have a better performance, due to its packet-processing optimized design and the circumvention of the kernel space.

Comparing the IPv6 results with IPv4, the performance in the case of IPv4 is slightly superior for different packets sizes (i.e., 64 and 128 bytes). This conclusion follows from the fact that the Key size in IPv6 is 128 bits increases the complexity of the LPM tree and affects the performance. Broadly translating, our findings in the performance tests

indicate that the CPU cycles consumed by or ODP LPM implementation are less than the DPDK implementation even in the case of IPv6, this is confirmed by the performance results where MACSAD presents a better result.

These experiments add to a growing corpus of research showing that P4 and ODP can work in conjunction to create a good base for a programmable dataplane device. P4 is gaining interest by the community, and it is in constant development and maintenance, as well as ODP that continues emerging as a good set of APIs for the dataplane backend.

The BB-Gen packet generator tool allows creating the PCAP and table trace files for the performance evaluation. BB-Gen can also create multiple sets of trace files with different packet sizes and different flow distribution simultaneously along with the corresponding table trace file for each PCAP and pipeline. BB-Gen effectively complements other tools of the P4 developers gadgets, such as `p4pktgen` that validates a P4 program by generating test input traffic and fills the tables, `p4app` that performs functional tests using BMV2 simple_switch in Mininet environments for P4 programs, and `P4benchmark` that tests the target compilers generating different P4 applications varying the complexity.

We tested the applicability of BB-Gen in real projects i.e., MACSAD and T4P4S, by creating the traces and PCAPs from a P4 file for different use cases with varying complexities, suggesting that other developers could benefit from our tool in their evaluation workflows, as well as the growing community striving for research reproducibility, e.g., by re-using traces, toolchains, and evaluation methodologies. With BB-Gen, NFPA and "MACSAD & T4P4S" we demonstrated that they can work conjointly to illustrate the datapath pipeline defined by a P4 code, create match+action table traces to control the P4 tables in the DUT, and generate a suitable PCAP file used for performance evaluation of different setups.

Altogether, this work accomplished a series of open source contributions: (1) the IPv6 lookup library developed will be suggested for adoption by ODP project, (2) BB-Gen, (3) P4 pipelines for IPv4 and IPv6, and (4) various trace files for the NFPA repository.

## Future Work

To conclude, we list some future tasks and open questions worth further investigation.

Related to BB-Gen:

- Integration with `P4Runtime` [1] to control the DUT, leveraging the possibility to configure the P4 table entries by adding, updating, deleting, along with others.

---

[1]   <https://p4.org/p4-runtime/>

- BB-Gen will allow to rely on Scapy's extensibility features to create packets with arbitrary sequence of headers that might not conform to standards, but an unorthodox P4 program may require it.

- Although Scapy is a powerful Python library to easily generate any type of packet, we have experienced that generating hundreds of thousands of packets can take several minutes, however since the performance evaluations are usually done offline (i.e., before an actual deployment), the measurements do not require real-time traffic generation. Therefore, we also take into consideration of Hexdump generation of packets proven to be more effective[2] and contribute to materialize other open-source projects into BB-gen.

Related to the LPM implementation:

- A deep analysis in how the performance is affected by the variation of the prefix length. In this work, we covered a brief overview of the impact of the prefix length in the performance. This provides a good starting point for discussion and further research, analyzing real distributions and how distribution can be optimized in a forwarding device.

- We defined the prefix tree levels following the ODP IPV4 specification. It will be important that future research investigate if there is an optimal level distribution for IPv4 and IPv6, and how this distribution can affect the performance of the device.

- With the evaluated performance properties, e.g., packet loss, latency, and CPU cycles. we demonstrated that the ODP implementation has good performance compared with other projects. Future research should further develop and confirm these initial findings by using the information to optimize even more the library.

- We tested MACSAD in an X86 environment; future research should evaluate in different platforms as ARMv8, Tofino, this can open a new area with the inclusion of more CPU cores and a dedicate P4 chipset.

- To improve an optimize the LPM implementation, the analysis of the Cache-misses and CPU cycles consumed by the process is an important part of the research. We evaluated our implementation using Perf; future research should explore other tools such as eBPF that can prepare user information in kernel context and transfer only needed information to user space adding more resources and statistics maps.

---

[2] <https://github.com/cslev/pcap_gen_perftest>

# Bibliography

BIONDI, P. *Welcome to Scapy's documentation!* 2008. Available from Internet: <http://scapy.readthedocs.io/en/latest/>. Cited 2 times on pages 18 and 35.

BOSSHART et al. P4: Programming protocol-independent packet processors. *SIGCOMM Comput. Commun. Rev.*, ACM, New York, NY, USA, jul. 2014. ISSN 0146-4833. Cited 2 times on pages 17 and 23.

BOTTA, A.; DAINOTTI, A.; PESCAPÉ, A. A tool for the generation of realistic network workload for emerging networking scenarios. *Computer Networks*, Elsevier, v. 56, n. 15, p. 3531–3547, 2012. Cited on page 67.

BRADNER, S.; MCQUAID, J. *Benchmarking Methodology for Network Interconnect Devices.* 1999. RFC 2544. Cited 4 times on pages 28, 68, 71, and 74.

CESEN, F. R.; PATRA, P. G. K.; ROTHENBERG, C. E. Bb-gen: A packet crafter for data plane evaluation. In: SBRC. [S.l.], 2018. Cited on page 67.

CSIKOR, L.; SZALAY, M.; SONKOLY, B.; TOKA, L. *Network Function Performance Analyzer.* 2015. <http://nfpa.tmit.bme.hu>. Cited 3 times on pages 28, 29, and 67.

CSIKOR, L.; SZALAY, M.; SONKOLY, B.; TOKA, L. Nfpa: Network function performance analyzer. *IEEE Conference on Network Function Virtualization and Software Defined Networks Demo Track*, 2015. Cited 4 times on pages 9, 28, 30, and 67.

DANG, H. T.; WANG, H.; JEPSEN, T.; BREBNER, G.; KIM, C.; REXFORD, J.; SOULÉ, R.; WEATHERSPOON, H. Whippersnapper: A p4 language benchmark suite. In: *Proceedings of the Symposium on SDN Research.* New York, NY, USA: ACM, 2017. (SOSR '17), p. 95–101. ISBN 978-1-4503-4947-5. Available from Internet: <http://doi.acm.org/10.1145/3050220.3050231>. Cited on page 36.

DEGERMARK, M.; BRODNIK, A.; CARLSSON, S.; PIN, S. Small forwarding tables for fast routing lookups. *ACM SIGCOMM Computer Communication Review*, 1997. Available from Internet: <http://dl.acm.org/citation.cfm?id=263133>. Cited on page 35.

EATHERTON GEORGE VARGHESE, Z. D. W. Tree bitmap : Hardware/software ip lookups with incremental updates. *ACM SIGCOMM Computer Communications Review*, April 2004. Available from Internet: <http://cseweb.ucsd.edu/~varghese/PAPERS/ccr2004.pdf>. Cited on page 34.

EATHERTONK, W.; VARGHESE, G.; DITTIA, Z. Tree bitmap: hardware/software ip lookups with incremental updates. *ACM SIGCOMM Computer Communications Review*, 2004. Available from Internet: <http://cseweb.ucsd.edu/~varghese/PAPERS/ccr2004.pdf>. Cited on page 35.

FERRAGINA, P.; LUCCIO, F.; MANZINI, G.; MUTHUKRISHNAN, S. Compressing and indexing labeled trees, with applications. *Journal of the ACM (JACM)*, November 2009. Available from Internet: <http://dl.acm.org/citation.cfm?id=1613680>. Cited on page 33.

FIESSLER, A.; HAGER, S.; SCHEUERMANN, B. Flexible line speed network packet classification using hybrid on-chip matching circuits. In: *HPSR '17: Proceedings of 2017 IEEE 18th International Conference on High Performance Switching and Routing.* [S.l.: s.n.], 2017. To appear.   Cited on page 35.

GNUPLOT. *Gnuplot.* 1986. <http://www.gnuplot.info/>.   Cited on page 29.

GUREVICH, V. *P4 Tutorial.* 2015. <https://p4.org/assets/ Nov-2015-P4-Bootcamp-P4-Tutorial.pdf>. Available from Internet: <https: //p4.org/assets/Nov-2015-P4-Bootcamp-P4-Tutorial.pdf>.   Cited 3 times on pages 9, 22, and 23.

HANNA, M.; CHO, S.; MELHEM, R. A novel scalable ipv6 lookup scheme using compressed pipelined tries. *IFIP 10th international TC 6 conference on Networking - Volume Part I*, May 2011. Available from Internet: <http: //dl.acm.org/citation.cfm?id=2008820>.   Cited on page 34.

INTEL. *Intel Threading Building Blocks.* 2011. <http://threadingbuildingblocks.org/>. Available from Internet: <http://threadingbuildingblocks.org/>.   Cited on page 34.

INTEL. *DPDK: Data Plane Development Kit.* 2014. <http://dpdk.org>.   Cited 2 times on pages 28 and 29.

JUNIPER NETWORKS. *How many Packets per Second per port are needed to achieve Wire-Speed?* 2015. Https://kb.juniper.net. Available from Internet: <https://kb.juniper.net/InfoCenter/index?page=content&id=KB14737>.   Cited on page 29.

KNUTSSON, K. *RWS Synthetic Pcap Generator.* 2014. Https://github.com/karknu/rws. Available from Internet: <https://github.com/karknu/rws>.   Cited on page 35.

KREUTZ, D.; RAMOS, F. M. V.; VERISSIMO, P.; ROTHENBERG, C. E.; AZODOLMOLKY, S.; UHLIG, S. Software-Defined Networking: A Comprehensive Survey. p. 1–61, 2014. Available from Internet: <http://arxiv.org/abs/1406.0440>. Cited on page 17.

LAKI, S.; HORPáCSI, D.; VöRöS, P.; KITLEI, R.; LESKó, D.; TEJFEL, M. High speed packet forwarding compiled from protocol independent data plane specifications. In: *ACM SIGCOMM'16 Posters and Demos.* [S.l.: s.n.], 2016.   Cited on page 35.

LETTIERI, G.; MAFFIONE, V.; RIZZO, L. A survey of fast packet i/o technologies for network function virtualization. In: KUNKEL, J. M.; YOKOTA, R.; TAUFER, M.; SHALF, J. (Ed.). *High Performance Computing.* Cham: Springer International Publishing, 2017. p. 579–590. ISBN 978-3-319-67630-2.   Cited on page 76.

MAHALINGAM, M.; DUTT, D.; DUDA, K.; AGARWAL, P.; KREEGER, L.; SRIDHAR, T.; BURSELL, M.; WRIGHT, C. *Virtual eXtensible Local Area Network (VXLAN): A Framework for Overlaying Virtualized Layer 2 Networks over Layer 3 Networks.* [S.l.], 2014. <https://tools.ietf.org/html/rfc7348>. Available from Internet: <https://tools.ietf.org/html/rfc7348>.   Cited on page 29.

MCCANNE, S.; JACOBSON, V. The bsd packet filter: A new architecture for user-level packet capture. In: *USENIX winter.* [S.l.: s.n.], 1993. v. 46.   Cited on page 80.

MCKEOWN, N. *Programming the Forwarding Plane.* 2016.
Https://forum.stanford.edu/events/2016/slides/plenary/Nick.pdf. Available from
Internet: <https://forum.stanford.edu/events/2016/slides/plenary/Nick.pdf>. Cited 3
times on pages 9, 22, and 25.

MELO, A. C. D. The new linux'perf'tools. In: *Slides from Linux Kongress.* [S.l.: s.n.],
2010. v. 18. Cited on page 80.

NIKOLENKO, S. I.; KOGANY, K.; RETVARIZ, G.; BERCZI-KOVACS, E. R.;
SHALIMOV, A. How to represent ipv6 forwarding tables on ipv4 or mpls dataplanes.
*IEEE Computer Communications Workshops (INFOCOM WKSHPS)*, April 2016.
Available from Internet: <http://ieeexplore.ieee.org/document/7562132/>. Cited on
page 34.

NÖTZLI, A.; KHAN, J.; FINGERHUT, A.; BARRETT, C.; ATHANAS, P. P4pktgen:
Automated test case generation for p4 programs. In: ACM. *Proceedings of the Symposium
on SDN Research.* [S.l.], 2018. p. 5. Cited on page 36.

OLSSON, R. Pktgen the linux packet generator. In: *Proceedings of the Linux Symposium,
Ottawa, Canada.* [S.l.: s.n.], 2005. v. 2, p. 11–24. Cited on page 67.

OPENDATAPLANE. *ODP Helper library.* 2013. <https://github.com/Linaro/odp/
tree/master/helper>. Available from Internet: <https://github.com/Linaro/odp/tree/
master/helper>. Cited on page 27.

OPENDATAPLANE. *OpenDataPlane DPDK.* 2013. Https://github.com/Linaro/odp-
dpdk. Available from Internet: <https://github.com/Linaro/odp-dpdk>. Cited on page
27.

OPENDATAPLANE. *OpenDataPlane.org.* 2013. <https://www.opendataplane.org>.
Cited 3 times on pages 9, 17, and 27.

OSTINATO. *Ostinato.* 2010. Https://ostinato.org. Available from Internet: <https:
//ostinato.org>. Cited on page 35.

OVS. *Open vSwitch.* 2009. <http://openvswitch.org/>. Available from Internet:
<http://openvswitch.org/>. Cited on page 33.

P4.ORG. *P4app.* 2013. Https://github.com/p4lang/p4app. Available from Internet:
<https://github.com/p4lang/p4app>. Cited on page 36.

PATRA, P. G.; ROTHENBERG, C. E.; PONGRáCZ, G. MACSAD: Multi-Architecture
Compiler System for Abstract Dataplanes (Aka Partnering P4 with ODP). In: *ACM
SIGCOMM'16 Demo and Poster Session.* [S.l.: s.n.], 2016. ISBN 978-1-4503-4193-6.
Cited 4 times on pages 9, 17, 18, and 31.

PATRA, P. G.; ROTHENBERG, C. E.; PONGRACZ, G. Macsad: High performance
dataplane applications on the move. In: *IEEE HPSR.* [S.l.: s.n.], 2017. p. 1–6. Cited 2
times on pages 17 and 18.

PATRA, P. G. K.; CESEN, F. E. R.; MEJIA, J. S.; FEFERMAN, D. L.; CSIKOR,
L.; ROTHENBERG, C. E.; PONGRACZ, G. Toward a sweet spot of data plane
programmability, portability, and performance: On the scalability of multi-architecture

p4 pipelines. *IEEE Journal on Selected Areas in Communications*, v. 36, n. 12, p. 2603–2611, Dec 2018. ISSN 0733-8716.  Cited on page 18.

PONG, F.; TZENG, N. Concise lookup tables for ipv4 and ipv6 longest prefix matching in scalable routers. *IEEE/ACM Transactions on Networking*, v. 20, n. 3, p. 729–741, June 2012. ISSN 1063-6692.  Cited 2 times on pages 18 and 34.

RAVIKUMAR, V.; MAHAPATRA, R. Tcam architecture for ip lookup using prefix properties. *IEEE*, August 2004. Available from Internet: <http://ieeexplore.ieee.org/abstract/document/1289292/>.  Cited on page 33.

RIPE NCC. *Visibility of Prefix Lengths in IPv4 and IPv6*. 2010. Https://labs.ripe.net. Available from Internet: <https://labs.ripe.net/Members/dbayer/visibility-of-prefix-lengths>.  Cited on page 38.

RIZZO, L. netmap: A novel framework for fast packet i/o. In: *USENIX ATC 12*. [S.l.: s.n.], 2012. ISBN 978-931971-93-5.  Cited on page 28.

RODRIGUEZ, F.; PATRA, P. G. K.; CSIKOR, L.; ROTHENBERG, C.; LAKI, P. V. S.; PONGRáCZ, G. Bb-gen: A packet crafter for p4 target evaluation. In: *Proceedings of the ACM SIGCOMM 2018 Conference on Posters and Demos*. New York, NY, USA: ACM, 2018. (SIGCOMM '18), p. 111–113. ISBN 978-1-4503-5915-3. Available from Internet: <http://doi.acm.org/10.1145/3234200.3234229>.  Cited on page 67.

ROSEN, E.; VISWANATHAN, A.; CALLON, R. *Multiprotocol Label Switching Architecture*. [S.l.], 2001. <https://tools.ietf.org/html/rfc3031>. Available from Internet: <https://tools.ietf.org/html/rfc3031>.  Cited on page 29.

RéTVáRI, G.; TAPOLCAI, J.; KORöSI, A.; MAJDáN, A.; HESZBERGER, Z. Compressing ip forwarding tables: Towards entropy bounds and beyond. *ACM SIGCOMM Computer Communication Review*, August 2013. Available from Internet: <http://conferences.sigcomm.org/sigcomm/2013/papers/sigcomm/p111.pdf>.  Cited 3 times on pages 18, 33, and 34.

SHAHHAZ et al. PISCES: A Programmable, Protocol-Independent Software Switch. In: *ACM SIGCOMM*. [S.l.: s.n.], 2016. ISBN 978-1-4503-4193-6.  Cited 3 times on pages 18, 25, and 33.

SKLOWER, K. *A tree-based packet routing table for Berkeley UNIX*. [S.l.], 2001. Available from Internet: <https://people.eecs.berkeley.edu/~sklower/routing.pdf>.  Cited on page 33.

SRINIVASAN, V.; VARGHESE, G. Faster ip lookups using controlled prefix expansion. *ACM SIGMETRICS Performance Evaluation Review*, 1999. Available from Internet: <http://dl.acm.org/citation.cfm?id=277863>.  Cited on page 34.

TCPDUMP. *Tcpdump*. 2010. Https://www.tcpdump.org. Available from Internet: <https://www.tcpdump.org>.  Cited on page 67.

TOBOLA, J.; KOřENEK, J. Effective hash-based ipv6 longest prefix match. *IEEE 14th International Symposium Design and Diagnostics of Electronic Circuits and Systems*, May 2011. Available from Internet: <http://ieeexplore.ieee.org/document/5783105/>.  Cited on page 34.

TURULL, D.; SJöDIN, P.; OLSSON, R. Pktgen: Measuring performance on high speed networks. *Computer Communications*, 2016. ISSN 0140-3664.  Cited on page 29.

WANG, G.; TZENG, N. feng. Tcam-based forwarding engine with minimum independent prefix set (mips) for fast updating. *Communications, 2006. ICC '06. IEEE International Conference on*, December 2006. Available from Internet: <http://ieeexplore.ieee.org/document/4024102/>.  Cited on page 35.

WATERLOO, U. of. *Binary trees*. 2018. <https://www.student.cs.uwaterloo. ca/~cs115/coursenotes1/08-bintrees-post.pdf>. Available from Internet: <https: //www.student.cs.uwaterloo.ca/~cs115/coursenotes1/08-bintrees-post.pdf>.  Cited 2 times on pages 9 and 39.

WILES, K. *pktgen-dpdk*. 2010. <http://dpdk.org/browse/apps/pktgen-dpdk/refs/>. Cited on page 29.

ZHOU, D.; FAN, B.; LIM, H.; KAMINSKY, M.; ANDERSEN, D. G. Scalable, high performance ethernet forwarding with cuckooswitch. *Proceedings of the ninth ACM conference on Emerging networking experiments and technologies*, December 2013. Available from Internet: <http://dl.acm.org/citation.cfm?id=2535379>.  Cited on page 33.

# Annex

# ANNEX A – Publications

- **F. R. Cesen**, G. P. Patra, and C. E. Rothenberg. BB-Gen: A Packet Crafter for Data Plane Evaluation. In: *XXXVI Simpósio Brasileiros de Redes de Computadores SBRC 2018 Salão de Ferramentas*, Campos do Jordão, SP, Brazil, May 2018.

- **F. R. Cesen**, G. P. Patra, C. E. Rothenberg, and G. Pongrácz. BB-Gen: A Packet Crafter for Performance Evaluation of P4 Data Planes. In: *5th P4 Workshop*, Stanford University, CA, USA, June 2018.

- G. P. patra, **F. R. Cesen**, J. S. Mejia, D. Feferman, C. E. Rothenberg, and G. Pongrácz. MACSAD: An Exemplar Realization of Multi-Architecture P4 Pipelines. In: *5th P4 Workshop*, Stanford University, CA, USA, June 2018.

- **F. R. Cesen**, G. P. Patra, C. E. Rothenberg, and G. Pongrácz. Design, Implementation and Evaluation of IPv4/IPv6 Longest Prefix Match support in P4 Dataplanes. In: *17º Workshop em Desempenho de Sistemas Computacionais e de Comunicação WPerformance 2018*, Natal, RN, Brazil, July 2018.

- **F. R. Cesen**, G. P. Patra, L. Csikor, C. E. Rothenberg, P. Vörös, S. Laki and G. Pongrácz. BB-Gen: A Packet Crafter for P4 Target Evaluation. In: *ACM Special Interest Group on Data Communication SIGCOMM 2018 Posters, Demos, and Student Research Competition*, Budapest, Hungary, August 2018.

- G. P. Patra, **F. R. Cesen**, J. S. Mejia, D. Feferman, L. Csikor, C. E. Rothenberg, and G. Pongrácz. Towards a Sweet Spot of Dataplane Programmability, Portability and Performance: On the Scalability of Multi-Architecture P4 Pipelines. Under Submission in: *IEEE COMSOC JSAC'18 Special Issue on Scalability Issues and Solutions for Software Defined Networks.*

- **F. E. R. Cesen**, and C. R. E. Rothenberg. Design, Implementation and Evaluation of IPv4/IPv6 Longest Prefix Match support in Multi-Architecture Programmable Dataplanes. In: *Décimo Primeiro Encontro dos Alunos e Docentes do Departamento de Engenharia de Computação e Automação Industrial*, Campinas, SP, Brazil, November 2018.

# ANNEX  B  −  ODP IPv6 Code

## B.1   ODP IPv6 lookup code

```
1   #include "config.h"
2   #include <string.h>
3   #include <stdint.h>
4   #include <errno.h>
5   #include <stdio.h>
6   #include <odp/helper/odph_iplookuptablev6.h>
7   #include "odph_list_internal.h"
8   #include "odph_debug.h"
9   #include <odp_api.h>
10  #include <odp/helper/ip.h>
11
12  typedef __int128 _uint128_t;
13  typedef unsigned __int128 uint128_t;
14
15  /** @magic word, write to the first byte of the memory block
16   *    to indicate this block is used by a ip lookup table
17   */
18  #define ODPH_IP_LOOKUP_TABLE_MAGIC_WORD 0xCFCFFCFC
19
20  /* The length(bit) of the IPv6 address */
21  #define IP_LENGTH 128
22
23  /* The number of L1 entries */
24  #define ENTRY_NUM_L1     (1 << 24)
25  /* The size of one L2\L3 subtree */
26  #define ENTRY_NUM_SUBTREE (1 << 12)
27
28  #define WHICH_CHILD(ip, cidr) ((ip >> (IP_LENGTH − cidr)) & 0x00000001)
29
30  /** @internal entry struct
31   *    Structure store an entry of the ip prefix table.
32   *    Because of the leaf pushing, each entry of the table must have
33   *    either a child entry, or a nexthop info.
34   *    If child == 0 and index != ODP_BUFFER_INVALID, this entry has
35   *    a nexthop info, index indicates the buffer that stores the
36   *    nexthop value, and ptr points to the address of the buffer.
37   *    If child == 1, this entry has a subtree, index indicates
38   *    the buffer that stores the subtree, and ptr points to the
39   *    address of the buffer.
40   */
```

```
41   typedef struct {
42     union {
43       odp_buffer_t nexthop;
44       void *ptr;
45     };
46     union {
47       _uint128_t u8;
48       struct {
49 #if ODP_BYTE_ORDER == ODP_BIG_ENDIAN
50         uint8_t child : 1;
51         uint8_t cidr  : 7;
52 #else
53         uint8_t cidr  : 7;
54         uint8_t child : 1;
55 #endif
56       };
57     };
58   } prefix_entry_t;
59
60   #define ENTRY_SIZE (sizeof(prefix_entry_t) + sizeof(odp_buffer_t))
61   #define ENTRY_BUFF_ARR(x) ((odp_buffer_t *)(void *)((char *)x \
62         + sizeof(prefix_entry_t) * ENTRY_NUM_SUBTREE))
63
64   /** @internal trie node struct
65    *   In this IP lookup algorithm, we use a
66    *   binary tire to detect the overlap prefix.
67    */
68   typedef struct trie_node {
69     /* tree structure */
70     struct trie_node *parent;
71     struct trie_node *left;
72     struct trie_node *right;
73     /* IP prefix length */
74     uint8_t cidr;
75     /* Nexthop buffer index */
76     odp_buffer_t nexthop;
77     /* Buffer that stores this node */
78     odp_buffer_t buffer;
79   } trie_node_t;
80
81   /** Number of L2\L3 entries(subtrees) per cache cube. */
82   #define CACHE_NUM_SUBTREE (1 << 13)
83   /** Number of trie nodes per cache cube. */
84   #define CACHE_NUM_TRIE    (1 << 20)
85
86   /** @typedef cache_type_t
87    *   Cache node type
```

```
 88     */
 89    typedef enum {
 90      CACHE_TYPE_SUBTREE = 0,
 91      CACHE_TYPE_TRIE
 92    } cache_type_t;
 93
 94    /** A IP lookup table structure. */
 95    typedef struct {
 96      /**< for check */
 97      uint32_t magicword;
 98      /** Name of the hash. */
 99      char name[ODPH_TABLE_NAME_LEN];
100      /** Total L1 entries. */
101      prefix_entry_t *l1e;
102      /** Root node of the binary trie */
103      trie_node_t *trie;
104      /** Length of value. */
105      uint32_t nexthop_len;
106      /** Queues of free slots (caches)
107       *   There are two queues:
108       *   - free_slots[CACHE_TYPE_SUBTREE] is used for L2 and
109       *     L3 entries (subtrees). Each entry stores an 8-bit
110       *     subtree.
111       *   - free_slots[CACHE_TYPE_TRIE] is used for the binary
112       *     trie. Each entry contains a trie node.
113       */
114      odp_queue_t free_slots[2];
115      /** The number of pool used by each queue. */
116      uint32_t cache_count[2];
117    } odph_iplookup_table_impl ODP_ALIGNED_CACHE;
118
119    /*************************************************************
120     ****************    Cache management    ********************
121     *************************************************************/
122
123    /** Destroy all caches */
124    static void cache_destroy(odph_iplookup_table_impl *impl){
125      odp_queue_t queue;
126      odp_event_t ev;
127      uint32_t i = 0, count = 0;
128      char pool_name[ODPH_TABLE_NAME_LEN + 8];
129
130      /* free all buffers in the queue */
131      for (; i < 2; i++) {
132        queue = impl->free_slots[i];
133        if (queue == ODP_QUEUE_INVALID)
134          continue;
```

```
135
136        while  (( ev = odp_queue_deq ( queue ) )
137            != ODP_EVENT_INVALID)  {
138          odp_buffer_free ( odp_buffer_from_event ( ev ) ) ;
139        }
140        odp_queue_destroy ( queue ) ;
141      }
142
143      /* destroy all cache pools */
144      for  ( i = 0;  i < 2;  i++) {
145        for  ( count = 0;  count < impl−>cache_count [ i ] ;  count++) {
146          sprintf (
147              pool_name ,  "%s_%d_%d" ,
148              impl−>name ,  i ,  count ) ;
149          odp_pool_destroy ( odp_pool_lookup ( pool_name ) ) ;
150        }
151      }
152    }
153
154    /** According to the type of cahce , set the value of
155     *  a buffer to the initial value .
156     */
157    static void cache_init_buffer ( odp_buffer_t buffer , cache_type_t type ,
158        uint32_t size ) {
158      int i = 0;
159      void ∗addr = odp_buffer_addr ( buffer ) ;
160
161      memset ( addr ,  0 ,  size ) ;
162      if  ( type == CACHE_TYPE_SUBTREE)  {
163        prefix_entry_t ∗entry = ( prefix_entry_t ∗) addr ;
164
165        for  ( i = 0;  i < ENTRY_NUM_SUBTREE;  i++, entry++)
166          entry−>nexthop = ODP_BUFFER_INVALID;
167      } else if  ( type == CACHE_TYPE_TRIE)  {
168        trie_node_t ∗node = ( trie_node_t ∗) addr ;
169
170        node−>buffer = buffer ;
171        node−>nexthop = ODP_BUFFER_INVALID;
172      }
173    }
174
175    /** Create a new buffer pool , and insert its buffer into the queue . */
176    static int cache_alloc_new_pool ( odph_iplookup_table_impl ∗tbl , cache_type_t
177        type ) {
177      odp_pool_t pool ;
178      odp_pool_param_t param ;
179      odp_queue_t queue = tbl−>free_slots [ type ] ;
```

```
180
181     odp_buffer_t buffer;
182     char pool_name[ODPH_TABLE_NAME_LEN + 8];
183     uint32_t size = 0, num = 0;
184
185     /* Create new pool (new free buffers). */
186     odp_pool_param_init(&param);
187     param.type = ODP_POOL_BUFFER;
188     param.buf.align = ODP_CACHE_LINE_SIZE;
189     if (type == CACHE_TYPE_SUBTREE) {
190       num = CACHE_NUM_SUBTREE;
191       size = ENTRY_SIZE * ENTRY_NUM_SUBTREE;
192     } else if (type == CACHE_TYPE_TRIE) {
193       num = CACHE_NUM_TRIE;
194       size = sizeof(trie_node_t);
195     } else {
196       ODPH_DBG("wrong cache_type_t.\n");
197       return -1;
198     }
199     param.buf.size = size;
200     param.buf.num = num;
201
202     sprintf(
203         pool_name, "%s_%d_%d",
204         tbl->name, type, tbl->cache_count[type]);
205     pool = odp_pool_create(pool_name, &param);
206     if (pool == ODP_POOL_INVALID) {
207       ODPH_DBG("failed to create a new pool.\n");
208       return -1;
209     }
210
211     /* insert new free buffers into queue */
212     while ((buffer = odp_buffer_alloc(pool))
213         != ODP_BUFFER_INVALID) {
214       cache_init_buffer(buffer, type, size);
215       odp_queue_enq(queue, odp_buffer_to_event(buffer));
216     }
217
218     tbl->cache_count[type]++;
219     return 0;
220 }
221
222 /** Get a new buffer from a cache list. If there is no
223  *   available buffer, allocate a new pool.
224  */
225 static odp_buffer_t cache_get_buffer(odph_iplookup_table_impl *tbl,
        cache_type_t type){
```

```
226      odp_buffer_t buffer = ODP_BUFFER_INVALID;
227      odp_queue_t queue = tbl->free_slots[type];
228
229      /* get free buffer from queue */
230      buffer = odp_buffer_from_event(
231            odp_queue_deq(queue));
232
233      /* If there is no free buffer available, allocate new pool */
234      if (buffer == ODP_BUFFER_INVALID) {
235        cache_alloc_new_pool(tbl, type);
236        buffer = odp_buffer_from_event(odp_queue_deq(queue));
237      }
238
239      return buffer;
240    }
241
242    /*************************************************************
243     ******************     Binary trie     ********************
244     *************************************************************/
245
246    /* Initialize the root node of the trie */
247    static int trie_init(odph_iplookup_table_impl *tbl){
248      trie_node_t *root = NULL;
249      odp_buffer_t buffer = cache_get_buffer(tbl, CACHE_TYPE_TRIE);
250
251      if (buffer != ODP_BUFFER_INVALID) {
252        root = (trie_node_t *)odp_buffer_addr(buffer);
253        root->cidr = 0;
254        tbl->trie = root;
255        return 0;
256      }
257
258      return -1;
259    }
260
261    /* Destroy the whole trie (recursively) */
262    static void trie_destroy(odph_iplookup_table_impl *tbl, trie_node_t *trie){
263      if (trie->left != NULL)
264        trie_destroy(tbl, trie->left);
265      if (trie->right != NULL)
266        trie_destroy(tbl, trie->right);
267
268      /* destroy this node */
269      odp_queue_enq(
270            tbl->free_slots[CACHE_TYPE_TRIE],
271            odp_buffer_to_event(trie->buffer));
272    }
```

```
273
274   /* Insert a new prefix node into the trie
275    * If the node is already existed, update its nexthop info,
276    *    Return 0 and set nexthop pointer to INVALID.
277    * If the node is not exitsed, create this target node and
278    *    all nodes along the path from root to the target node.
279    *    Then return 0 and set nexthop pointer points to the
280    *    new buffer.
281    * Return -1 for error.
282    */
283   static int trie_insert_node(
284       odph_iplookup_table_impl *tbl, trie_node_t *root,
285       _uint128_t ip, uint8_t cidr, odp_buffer_t nexthop)
286   {
287     uint8_t level = 0, child;
288     odp_buffer_t buf;
289     trie_node_t *node = root, *prev = root;
290
291     /* create/update all nodes along the path
292      * from root to the new node. */
293     for (level = 1; level <= cidr; level++) {
294       child = WHICH_CHILD(ip, level);
295
296       node = child == 0 ? prev->left : prev->right;
297       /* If the child node doesn't exit, create it. */
298       if (node == NULL) {
299         buf = cache_get_buffer(tbl, CACHE_TYPE_TRIE);
300         if (buf == ODP_BUFFER_INVALID)
301           return -1;
302
303         node = (trie_node_t *)odp_buffer_addr(buf);
304         node->cidr = level;
305         node->parent = prev;
306
307         if (child == 0)
308           prev->left = node;
309         else
310           prev->right = node;
311       }
312       prev = node;
313     }
314
315     /* The final one is the target. */
316     node->nexthop = nexthop;
317     return 0;
318   }
319
```

```
320    /* Delete a node */
321    static int trie_delete_node(
322        odph_iplookup_table_impl *tbl,
323        trie_node_t *root, _uint128_t ip, uint8_t cidr)
324    {
325      if (root == NULL)
326        return -1;
327
328      /* The default prefix (root node) cannot be deleted. */
329      if (cidr == 0)
330        return -1;
331
332      trie_node_t *node = root, *prev = NULL;
333      uint8_t level = 1, child = 0;
334      odp_buffer_t tmp;
335
336      /* Find the target node. */
337      for (level = 1; level <= cidr; level++) {
338        child = WHICH_CHILD(ip, level);
339        node = (child == 0) ? node->left : node->right;
340        if (node == NULL) {
341          ODPH_DBG("Trie node is not existed\n");
342          return -1;
343        }
344      }
345
346      node->nexthop = ODP_BUFFER_INVALID;
347
348      /* Delete all redundant nodes along the path. */
349      for (level = cidr; level > 0; level--) {
350        if (
351          node->left != NULL || node->right != NULL ||
352          node->nexthop != ODP_BUFFER_INVALID)
353          break;
354
355        child = WHICH_CHILD(ip, level);
356        prev = node->parent;
357
358        /* free trie node */
359        tmp = node->buffer;
360        cache_init_buffer(
361            tmp, CACHE_TYPE_TRIE, sizeof(trie_node_t));
362        odp_queue_enq(
363            tbl->free_slots[CACHE_TYPE_TRIE],
364            odp_buffer_to_event(tmp));
365
366        if (child == 0)
```

```
367          prev->left = NULL;
368        else
369          prev->right = NULL;
370        node = prev;
371      }
372      return 0;
373  }
374
375  /* Detect the longest overlapping prefix. */
376  static int trie_detect_overlap(
377      trie_node_t *trie, _uint128_t ip, uint8_t cidr,
378      uint8_t leaf_push, uint8_t *over_cidr,
379      odp_buffer_t *over_nexthop)
380  {
381    uint8_t child = 0;
382    uint32_t level, limit = cidr > leaf_push ? leaf_push + 1 : cidr;
383    trie_node_t *node = trie, *longest = trie;
384
385    for (level = 1; level < limit; level++) {
386      child = WHICH_CHILD(ip, level);
387      node = (child == 0) ? node->left : node->right;
388      if (node->nexthop != ODP_BUFFER_INVALID)
389        longest = node;
390    }
391
392    *over_cidr = longest->cidr;
393    *over_nexthop = longest->nexthop;
394    return 0;
395  }
396
397  /***********************************************************
398   ***************    IP prefix lookup table    ***************
399   ***********************************************************/
400
401  odph_table_t odph_iplookupv6_table_lookup(const char *name) {
402    odph_iplookup_table_impl *tbl = NULL;
403    odp_shm_t shm;
404
405    if (name == NULL || strlen(name) >= ODPH_TABLE_NAME_LEN)
406      return NULL;
407
408    shm = odp_shm_lookup(name);
409    if (shm != ODP_SHM_INVALID)
410      tbl = (odph_iplookup_table_impl *)odp_shm_addr(shm);
411
412    if (
413      tbl != NULL &&
```

```
414        tbl->magicword == ODPH_IP_LOOKUP_TABLE_MAGIC_WORD &&
415        strcmp(tbl->name, name) == 0)
416        return (odph_table_t)tbl;
417
418    return NULL;
419  }
420
421  odph_table_t odph_iplookupv6_table_create(const char *name,
422            uint32_t p1 ODP_UNUSED,
423            uint32_t p2 ODP_UNUSED,
424            uint32_t value_size)
425  {
426    odph_iplookup_table_impl *tbl;
427    odp_shm_t shm_tbl;
428    odp_queue_t queue;
429    odp_queue_param_t qparam;
430    unsigned i;
431    uint32_t impl_size, l1_size;
432    char queue_name[ODPH_TABLE_NAME_LEN + 2];
433
434    /* Check for valid parameters */
435    if (strlen(name) == 0) {
436      ODPH_DBG("invalid parameters\n");
437      return NULL;
438    }
439
440    /* Guarantee there's no existing */
441    tbl = (void *)odph_iplookupv6_table_lookup(name);
442    if (tbl != NULL) {
443      ODPH_DBG("IP prefix table %s already exists\n", name);
444      return NULL;
445    }
446
447    /* Calculate the sizes of different parts of IP prefix table */
448    impl_size = sizeof(odph_iplookup_table_impl);
449    l1_size = ENTRY_SIZE * ENTRY_NUM_L1;
450
451    shm_tbl = odp_shm_reserve(
452          name, impl_size + l1_size,
453          ODP_CACHE_LINE_SIZE, ODP_SHM_SW_ONLY);
454
455    if (shm_tbl == ODP_SHM_INVALID) {
456      ODPH_DBG(
457        "shm allocation failed for odph_iplookup_table_impl %s\n",
458        name);
459      return NULL;
460    }
```

```
461
462     tbl = (odph_iplookup_table_impl *)odp_shm_addr(shm_tbl);
463     memset(tbl, 0, impl_size + l1_size);
464
465     /* header of this mem block is the table impl struct,
466      * then the l1 entries array.
467      */
468     tbl->l1e = (prefix_entry_t *)(void *)((char *)tbl + impl_size);
469     for (i = 0; i < ENTRY_NUM_L1; i++)
470       tbl->l1e[i].nexthop = ODP_BUFFER_INVALID;
471
472     /* Setup table context. */
473     snprintf(tbl->name, sizeof(tbl->name), "%s", name);
474     tbl->magicword = ODPH_IP_LOOKUP_TABLE_MAGIC_WORD;
475     tbl->nexthop_len = value_size;
476
477     /* Initialize cache */
478     for (i = 0; i < 2; i++) {
479       tbl->cache_count[i] = 0;
480
481       odp_queue_param_init(&qparam);
482       qparam.type = ODP_QUEUE_TYPE_PLAIN;
483       sprintf(queue_name, "%s_%d", name, i);
484       queue = odp_queue_create(queue_name, &qparam);
485       if (queue == ODP_QUEUE_INVALID) {
486         ODPH_DBG("failed to create queue");
487         cache_destroy(tbl);
488         return NULL;
489       }
490       tbl->free_slots[i] = queue;
491       cache_alloc_new_pool(tbl, i);
492     }
493
494     /* Initialize tire */
495     if (trie_init(tbl) < 0) {
496       odp_shm_free(shm_tbl);
497       return NULL;
498     }
499
500     return (odph_table_t)tbl;
501   }
502
503   int odph_iplookupv6_table_destroy(odph_table_t tbl) {
504     int i, j;
505     odph_iplookup_table_impl *impl = NULL;
506     prefix_entry_t *subtree = NULL;
507     odp_buffer_t *buff1 = NULL, *buff2 = NULL;
```

```
508
509      if (tbl == NULL)
510        return −1;
511
512      impl = (odph_iplookup_table_impl *)(void *)tbl;
513
514      /* check magic word */
515      if (impl−>magicword != ODPH_IP_LOOKUP_TABLE_MAGIC_WORD) {
516        ODPH_DBG("wrong magicword for IP prefix table\n");
517        return −1;
518      }
519
520      /* destroy trie */
521      trie_destroy(impl, impl−>trie);
522
523      /* free all L2 and L3 entries */
524      buff1 = ENTRY_BUFF_ARR(impl−>l1e);
525      for (i = 0; i < ENTRY_NUM_L1; i++) {
526        if ((impl−>l1e[i]).child == 0)
527          continue;
528
529        subtree = (prefix_entry_t *)impl−>l1e[i].ptr;
530        buff2 = ENTRY_BUFF_ARR(subtree);
531        /* destroy all l3 subtrees of this l2 subtree */
532        for (j = 0; j < ENTRY_NUM_SUBTREE; j++) {
533          if (subtree[j].child == 0)
534            continue;
535          odp_queue_enq(
536              impl−>free_slots[CACHE_TYPE_TRIE],
537              odp_buffer_to_event(buff2[j]));
538        }
539        /* destroy this l2 subtree */
540        odp_queue_enq(
541            impl−>free_slots[CACHE_TYPE_TRIE],
542            odp_buffer_to_event(buff1[i]));
543      }
544
545      /* destroy all cache */
546      cache_destroy(impl);
547
548      /* free impl */
549      odp_shm_free(odp_shm_lookup(impl−>name));
550      return 0;
551    }
552
553    /* Insert the prefix into level x
554     * Return:
```

```
555    *   −1 error
556    *    0   the table is unmodified
557    *    1   the table is modified
558    */
559   static int prefix_insert_into_lx(
560       odph_iplookup_table_impl *tbl, prefix_entry_t *entry,
561       uint8_t cidr, odp_buffer_t nexthop, uint8_t level)
562   {
563     uint8_t ret = 0;
564     uint32_t i = 0, limit = (1 << (level − cidr));
565     prefix_entry_t *e = entry, *ne = NULL;
566
567     for (i = 0; i < limit; i++, e++) {
568       if (e−>child == 1) {
569         if (e−>cidr > cidr)
570           continue;
571
572         e−>cidr = cidr;
573         /* push to next level */
574         ne = (prefix_entry_t *)e−>ptr;
575         ret = prefix_insert_into_lx(
576             tbl, ne, cidr, nexthop, cidr + 8);
577       } else {
578         if (e−>cidr > cidr)
579           continue;
580
581         e−>child = 0;
582         e−>cidr = cidr;
583         e−>nexthop = nexthop;
584         ret = 1;
585       }
586     }
587     return ret;
588   }
589
590   static int prefix_insert_iter(
591       odph_iplookup_table_impl *tbl, prefix_entry_t *entry,
592       odp_buffer_t *buff, _uint128_t ip, uint8_t cidr,
593       odp_buffer_t nexthop, uint8_t level, uint8_t depth)
594   {
595     uint8_t state = 0;
596     prefix_entry_t *ne = NULL;
597     odp_buffer_t *nbuff = NULL;
598
599     /* If child subtree is existed, get it. */
600     if (entry−>child) {
601       //printf("child existed, get it. \n");
```

```
602        ne = (prefix_entry_t *)entry->ptr;
603        nbuff = ENTRY_BUFF_ARR(ne);
604      } else {
605        /* If the child is not existed, create a new subtree. */
606        odp_buffer_t buf, push = entry->nexthop;
607
608        buf = cache_get_buffer(tbl, CACHE_TYPE_SUBTREE);
609        if (buf == ODP_BUFFER_INVALID) {
610          ODPH_DBG("failed to get subtree buffer from cache.\n");
611          return -1;
612        }
613        ne = (prefix_entry_t *)odp_buffer_addr(buf);
614        nbuff = ENTRY_BUFF_ARR(ne);
615
616        entry->child = 1;
617        entry->ptr = ne;
618        *buff = buf;
619        /* If this entry contains a nexthop and a small cidr,
620         * push it to the next level.
621         */
622        if (entry->cidr > 0)
623          (void)prefix_insert_into_lx(tbl, ne, entry->cidr,
624                     push, entry->cidr + 8);
625      }
626      ne += (ip >> 120);
627      nbuff += (ip >> 120);
628      if (cidr <= 8) {
629        state = prefix_insert_into_lx(
630            tbl, ne, cidr + depth * 8, nexthop, level);
631      } else {
632        state = prefix_insert_iter(
633            tbl, ne, nbuff, ip << 8, cidr - 8,
634            nexthop, level + 8, depth + 1);
635      }
636
637      return state;
638    }
639
640    int odph_iplookupv6_table_put_value(odph_table_t tbl, void *key, void *
         value) {
641      odph_iplookup_table_impl *impl = (void *)tbl;
642      odph_iplookupv6_prefix_t *prefix = (odph_iplookupv6_prefix_t *)key;
643      prefix_entry_t *l1e = NULL;
644      odp_buffer_t nexthop;
645      int ret = 0;
646      _uint128_t lkp_ip = 0;
647
```

```
648      if ((tbl == NULL) || (key == NULL) || (value == NULL))
649        return -1;
650      nexthop = *((odp_buffer_t *)value);
651
652      if (prefix->cidr == 0)
653        return -1;
654
655      ret = odph_ipv6_addr_parse(&lkp_ip, "ffff:ffff:ffff:ffff:ffff:ffff:ffff:
          ffff");
656      if (ret < 0) {
657        printf("Failed to get IPv6 addr from str\n");
658        return -1;
659      }
660
661      prefix->ip = prefix->ip & (lkp_ip << (IP_LENGTH - prefix->cidr));
662      ret = trie_insert_node(
663            impl, impl->trie,
664            prefix->ip, prefix->cidr, nexthop);
665
666      if (ret < 0) {
667        ODPH_DBG("failed to insert into trie\n");
668        return -1;
669      }
670
671      l1e = &impl->l1e[prefix->ip >> 112];
672      odp_buffer_t *buff = ENTRY_BUFF_ARR(impl->l1e) + (prefix->ip >> 112);
673
674      if (prefix->cidr <= 16) {
675        ret = prefix_insert_into_lx(
676            impl, l1e, prefix->cidr, nexthop, 16);
677      } else {
678        ret = prefix_insert_iter(
679            impl, l1e, buff,
680            ((prefix->ip) << 112), prefix->cidr - 16,
681            nexthop, 24, 2);
682      }
683      return ret;
684  }
685
686  int odph_iplookupv6_table_get_value(odph_table_t tbl, void *key,
687            void *buffer ODP_UNUSED,
688            uint32_t buffer_size ODP_UNUSED)
689  {
690      odph_iplookup_table_impl *impl = (void *)tbl;
691      _uint128_t ip;
692      prefix_entry_t *entry;
693      odp_buffer_t *buff = (odp_buffer_t *)buffer;
```

```
694
695      if ((tbl == NULL) || (key == NULL) || (buffer == NULL))
696        return −EINVAL;
697      ip = *((_uint128_t *)key);
698      entry = &impl−>l1e[ip >> 112];
699      if (entry == NULL) {
700        ODPH_DBG("failed to get L1 entry.\n");
701        return −1;
702      }
703      ip <<= 112;
704
705      while (entry−>child) {
706        entry = (prefix_entry_t *)entry−>ptr;
707        entry += ip >> 24;
708        ip <<= 8;
709      }
710      if (entry−>nexthop == ODP_BUFFER_INVALID) {
711        /* ONLY match the default prefix */
712        printf("only match the default prefix\n");
713        *buff = ODP_BUFFER_INVALID;
714      } else {
715        *buff = entry−>nexthop;
716      }
717
718      return 0;
719    }
720
721    static int prefix_delete_lx(
722        odph_iplookup_table_impl *tbl, prefix_entry_t *l1e,
723        odp_buffer_t *buff, uint8_t cidr, uint8_t over_cidr,
724        odp_buffer_t over_nexthop, uint8_t level)
725    {
726      uint8_t ret, flag = 1;
727      prefix_entry_t *e = l1e;
728      odp_buffer_t *b = buff;
729      uint32_t i = 0, limit = 1 << (level − cidr);
730
731      for (i = 0; i < limit; i++, e++, b++) {
732        if (e−>child == 1) {
733          if (e−>cidr > cidr) {
734            flag = 0;
735            continue;
736          }
737
738          prefix_entry_t *ne = (prefix_entry_t *)e−>ptr;
739          odp_buffer_t *nbuff = ENTRY_BUFF_ARR(ne);
740
```

```
741        e->cidr = over_cidr;
742        ret = prefix_delete_lx(
743            tbl, ne, nbuff, cidr, over_cidr,
744            over_nexthop, cidr + 8);
745
746        /* If ret == 1, the next 2^8 entries equal to
747         * (over_cidr, over_nexthop). In this case, we
748         * should not push the (over_cidr, over_nexthop)
749         * to the next level. In fact, we should recycle
750         * the next 2^8 entries.
751         */
752        if (ret) {
753          /* destroy subtree */
754          cache_init_buffer(
755            *b, CACHE_TYPE_SUBTREE,
756            ENTRY_SIZE * ENTRY_NUM_SUBTREE);
757          odp_queue_enq(
758            tbl->free_slots[CACHE_TYPE_SUBTREE],
759            odp_buffer_to_event(*b));
760          e->child = 0;
761          e->nexthop = over_nexthop;
762        } else {
763          flag = 0;
764        }
765      } else {
766        if (e->cidr > cidr) {
767          flag = 0;
768          continue;
769        } else {
770          e->cidr = over_cidr;
771          e->nexthop = over_nexthop;
772        }
773      }
774    }
775    return flag;
776 }
777
778 /* Check if the entry can be recycled.
779  * An entry can be recycled duo to two reasons:
780  * - all children of the entry are the same,
781  * - all children of the entry have a cidr smaller than the level
782  *   bottom bound.
783  */
784 static uint8_t can_recycle(prefix_entry_t *e, uint32_t level) {
785    uint8_t recycle = 1;
786    int i = 1;
787    prefix_entry_t *ne = (prefix_entry_t *)e->ptr;
```

```
788
789    if (ne->child)
790      return 0;
791
792    uint8_t cidr = ne->cidr;
793    odp_buffer_t index = ne->nexthop;
794
795    if (cidr > level)
796      return 0;
797
798    ne++;
799    for (; i < 256; i++, ne++) {
800      if (
801          ne->child != 0 || ne->cidr != cidr ||
802          ne->nexthop != index) {
803        recycle = 0;
804        break;
805      }
806    }
807    return recycle;
808  }
809
810  static uint8_t prefix_delete_iter(
811      odph_iplookup_table_impl *tbl, prefix_entry_t *e,
812      odp_buffer_t *buff, _uint128_t ip, uint8_t cidr,
813      uint8_t level, uint8_t depth)
814  {
815    uint8_t ret = 0, over_cidr;
816    odp_buffer_t over_nexthop;
817
818    trie_detect_overlap(
819        tbl->trie, ip, cidr + 8 * depth, level,
820        &over_cidr, &over_nexthop);
821    if (cidr > 8) {
822      prefix_entry_t *ne =
823        (prefix_entry_t *)e->ptr;
824      odp_buffer_t *nbuff = ENTRY_BUFF_ARR(ne);
825
826      ne += ((_uint128_t)(ip << level) >> 24);
827      nbuff += ((_uint128_t)(ip << level) >> 24);
828      ret = prefix_delete_iter(
829          tbl, ne, nbuff, ip, cidr - 8,
830          level + 8, depth + 1);
831
832      if (ret && can_recycle(e, level)) {
833        cache_init_buffer(
834          *buff, CACHE_TYPE_SUBTREE,
```

```
835            ENTRY_SIZE * ENTRY_NUM_SUBTREE);
836        odp_queue_enq(
837          tbl->free_slots[CACHE_TYPE_SUBTREE],
838          odp_buffer_to_event(*buff));
839        e->child = 0;
840        e->nexthop = over_nexthop;
841        e->cidr = over_cidr;
842        return 1;
843      }
844      return 0;
845    }
846
847    ret = prefix_delete_lx(
848        tbl, e, buff, cidr + 8 * depth,
849        over_cidr, over_nexthop, level);
850    return ret;
851  }
852
853  int odph_iplookupv6_table_remove_value(odph_table_t tbl, void *key) {
854    odph_iplookup_table_impl *impl = (void *)tbl;
855    odph_iplookupv6_prefix_t *prefix = (odph_iplookupv6_prefix_t *)key;
856    _uint128_t ip;
857    uint8_t cidr;
858
859    if ((tbl == NULL) || (key == NULL))
860      return -EINVAL;
861
862    ip   = prefix->ip;
863    cidr = prefix->cidr;
864
865    if (cidr == 0)
866      return -EINVAL;
867
868    prefix_entry_t *entry = &impl->l1e[ip >> 112];
869    odp_buffer_t *buff = ENTRY_BUFF_ARR(impl->l1e) + (ip >> 112);
870    uint8_t over_cidr, ret;
871    odp_buffer_t over_nexthop;
872
873    trie_detect_overlap(
874        impl->trie, ip, cidr, 16, &over_cidr, &over_nexthop);
875
876    if (cidr <= 16) {
877      prefix_delete_lx(
878        impl, entry, buff, cidr, over_cidr, over_nexthop, 16);
879    } else {
880      prefix_entry_t *ne = (prefix_entry_t *)entry->ptr;
881      odp_buffer_t *nbuff = ENTRY_BUFF_ARR(ne);
```

```
882
883        ne += ((_uint128_t)(ip << 112) >> 24);
884        nbuff += ((_uint128_t)(ip << 112) >> 24);
885        ret = prefix_delete_iter(impl, ne, nbuff, ip, cidr − 16, 24, 2);
886
887        if (ret && can_recycle(entry, 16)) {
888          cache_init_buffer(
889            *buff, CACHE_TYPE_SUBTREE,
890            sizeof(prefix_entry_t) * ENTRY_NUM_SUBTREE);
891          odp_queue_enq(
892            impl−>free_slots[CACHE_TYPE_SUBTREE],
893            odp_buffer_to_event(*buff));
894          entry−>child = 0;
895          entry−>cidr = over_cidr;
896          entry−>nexthop = over_nexthop;
897        }
898      }
899
900      return trie_delete_node(impl, impl−>trie, ip, cidr);
901    }
902
903    odph_table_ops_t odph_iplookupv6_table_ops = {
904      odph_iplookupv6_table_create,
905      odph_iplookupv6_table_lookup,
906      odph_iplookupv6_table_destroy,
907      odph_iplookupv6_table_put_value,
908      odph_iplookupv6_table_get_value,
909      odph_iplookupv6_table_remove_value
910    };
```

Listing B.1 – ODP IPv6 lookup code

## B.2   ODP IPv6 Parse code

```
1   #include "config.h"
2   #include <odp/helper/ip.h>
3   #include <stdio.h>
4   #include <string.h>
5
6   typedef __int128 _uint128_t;
7   typedef unsigned __int128 uint128_t;
8
9   int odph_ipv4_addr_parse(uint32_t *ip_addr, const char *str) {
10    unsigned byte[ODPH_IPV4ADDR_LEN];
11    int i;
12
13    memset(byte, 0, sizeof(byte));
14
```

```
15    if (sscanf(str, "%u.%u.%u.%u",
16        &byte[0], &byte[1], &byte[2], &byte[3]) != ODPH_IPV4ADDR_LEN)
17      return -1;
18
19    for (i = 0; i < ODPH_IPV4ADDR_LEN; i++)
20      if (byte[i] > 255)
21        return -1;
22
23    *ip_addr = byte[0] << 24 | byte[1] << 16 | byte[2] << 8 | byte[3];
24
25    return 0;
26  }
27
28  int odph_ipv6_addr_parse(_uint128_t *ip_addr, const char *str){
29    unsigned byte[ODPH_IPV6ADDR_LEN];
30    int i;
31    _uint128_t p_ip1, p_ip2, p_ip3, p_ip4;
32
33    memset(byte, 0, sizeof(byte));
34
35    if (sscanf(str, "%02x%02x:%02x%02x:%02x%02x:%02x%02x:%02x%02x:%02x%02x
      :%02x%02x:%02x%02x",
36        &byte[0], &byte[1], &byte[2], &byte[3], &byte[4], &byte[5], &byte
      [6], &byte[7], &byte[8], &byte[9],
37        &byte[10], &byte[11], &byte[12], &byte[13], &byte[14], &byte[15]) !=
      ODPH_IPV6ADDR_LEN)
38      return -1;
39
40    for (i = 0; i < ODPH_IPV6ADDR_LEN; i++)
41      if (byte[i] > 255)
42        return -1;
43
44    p_ip1 = byte[0] << 24 | byte[1] << 16 | byte[2] << 8 | byte[3];
45    p_ip2 = byte[4] << 24 | byte[5] << 16 | byte[6] << 8 | byte[7];
46    p_ip3 = byte[8] << 24 | byte[9] << 16 | byte[10] << 8 | byte[11];
47    p_ip4 = byte[12] << 24 | byte[13] << 16 | byte[14] << 8 | byte[15];
48    *ip_addr = p_ip1 << 96 | p_ip2 << 64 | p_ip3 << 32 | p_ip4;
49
50    return 0;
51  }
```

Listing B.2 – ODP IPv6 parse code

# ANNEX C – Controller Code

## C.1   MACSAD IPv4 Controller code

```c
1  #include "controller.h"
2  #include "messages.h"
3  #include <unistd.h>
4  #include <stdio.h>
5  #include <string.h>
6  #include <time.h>
7
8  #define MAX_MACS 2000000
9
10 controller c;
11
12 void fill_ipv4_fib_lpm_table(uint8_t ip[4], uint8_t port, uint8_t mac[6])
13 {
14   char buffer[2048];
15   struct p4_header* h;
16   struct p4_add_table_entry* te;
17   struct p4_action* a;
18   struct p4_action_parameter* ap,* ap2;
19   struct p4_field_match_exact* exact;
20
21   h = create_p4_header(buffer, 0, 2048);
22   te = create_p4_add_table_entry(buffer,0,2048);
23   strcpy(te->table_name, "ipv4_fib_lpm");
24
25   exact = add_p4_field_match_exact(te, 2048);
26   strcpy(exact->header.name, "ipv4.dstAddr");
27   memcpy(exact->bitmap, ip, 4);
28   exact->length = 4*8+0;
29
30   a = add_p4_action(h, 2048);
31   strcpy(a->description.name, "fib_hit_nexthop");
32
33   ap = add_p4_action_parameter(h, a, 2048);
34   strcpy(ap->name, "dmac");
35   memcpy(ap->bitmap, mac, 6);
36   ap->length = 6*8+0;
37
38   ap2 = add_p4_action_parameter(h, a, 2048);
39   strcpy(ap2->name, "port");
40   ap2->bitmap[0] = port;
```

```
41    ap2->bitmap[1] = 0;
42    ap2->length = 2*8+0;
43
44    netconv_p4_header(h);
45    netconv_p4_add_table_entry(te);
46    netconv_p4_field_match_exact(exact);
47    netconv_p4_action(a);
48    netconv_p4_action_parameter(ap);
49    netconv_p4_action_parameter(ap2);
50
51    send_p4_msg(c, buffer, 2048);
52  }
53
54  void fill_sendout_table(uint8_t port, uint8_t smac[6])
55  {
56      char buffer[2048];
57      struct p4_header* h;
58      struct p4_add_table_entry* te;
59      struct p4_action* a;
60    struct p4_action_parameter* ap;
61    struct p4_field_match_exact* exact;
62
63      h = create_p4_header(buffer, 0, 2048);
64      te = create_p4_add_table_entry(buffer,0,2048);
65      strcpy(te->table_name, "sendout");
66
67      exact = add_p4_field_match_exact(te, 2048);
68      strcpy(exact->header.name, "standard_metadata.egress_port");
69      exact->bitmap[0] = port;
70    exact->bitmap[1] = 0;
71      exact->length = 2*8+0;
72
73    a = add_p4_action(h, 2048);
74    strcpy(a->description.name, "rewrite_src_mac");
75
76    ap = add_p4_action_parameter(h, a, 2048);
77    strcpy(ap->name, "smac");
78    memcpy(ap->bitmap, smac, 6);
79    ap->length = 6*8+0;
80
81      netconv_p4_header(h);
82      netconv_p4_add_table_entry(te);
83      netconv_p4_field_match_exact(exact);
84      netconv_p4_action(a);
85      netconv_p4_action_parameter(ap);
86      send_p4_msg(c, buffer, 2048);
87  }
```

```
88
89   uint8_t macs[MAX_MACS][6];
90   uint8_t portmap[MAX_MACS];
91   uint8_t ips[MAX_MACS][4];
92   int mac_count = -1;
93
94   int read_macs_and_ports_from_file(char *filename) {
95     FILE *f;
96     char line[200];
97     int values[6];
98     int values_ip[4];
99     int port;
100    int i;
101
102    f = fopen(filename, "r");
103    if (f == NULL) return -1;
104
105    while (fgets(line, sizeof(line), f)) {
106      line[strlen(line)-1] = '\0';
107      //TODO why %c?
108      if (11 == sscanf(line, "%d.%d.%d.%d %x:%x:%x:%x:%x:%x %d",
109            &values_ip[0], &values_ip[1], &values_ip[2], &values_ip[3],
110            &values[0], &values[1], &values[2],
111            &values[3], &values[4], &values[5], &port) )
112      {
113        if (mac_count==MAX_MACS-1)
114        {
115          printf("Too many entries...\n");
116          break;
117        }
118
119        ++mac_count;
120        for( i = 0; i < 6; ++i )
121          macs[mac_count][i] = (uint8_t) values[i];
122        for( i = 0; i < 4; ++i )
123          ips[mac_count][i] = (uint8_t) values_ip[i];
124        portmap[mac_count] = (uint8_t) port;
125
126      } else {
127        printf("Wrong format error in line %d : %s\n", mac_count+2, line);
128        fclose(f);
129        return -1;
130      }
131    }
132    fclose(f);
133    return 0;
134  }
```

```
135
136    void dhf(void* b) {
137      printf("Unknown digest received\n");
138    }
139
140    void init() {
141      int i;
142      uint8_t smac[6] = {0xd0, 0x69, 0x0f, 0xa8, 0x39, 0x90};
143      printf("INIT");
144      clock_t begin = clock();
145      for (i=0;i<=mac_count;++i)
146      {
147        fill_ipv4_fib_lpm_table(ips[i], portmap[i], macs[i]);
148
149              if (0 == (i%1000)){
150          printf("%d inside sleep \n",i);
151          sleep(1);
152        }
153        fill_sendout_table(portmap[i], smac);
154              usleep(1000);
155      }
156
157      clock_t end = clock();
158      double time_spent = (double)(end - begin) / CLOCKS_PER_SEC;
159      printf ("ctrl Total entries sent %d time %f\n",i,time_spent);
160
161    }
162
163    int main(int argc, char* argv[])
164    {
165      uint8_t ip[4] = {192,168,1,1};
166        uint8_t mac[6] = {0xa0, 0x36, 0x9f, 0x3e, 0x94, 0xea};
167      uint8_t port = 1;
168
169      uint8_t ip2[4] = {192,168,0,1};
170        uint8_t mac2[6] = {0xa0, 0x36, 0x9f, 0x3e, 0x94, 0xe8};
171      uint8_t port2 = 0;
172
173        uint8_t smac[6] = {0xd0, 0x69, 0x0f, 0xa8, 0x39, 0x90};
174
175
176      if (argc>1) {
177        if (argc!=2) {
178          printf("Too many arguments...\nUsage: %s <filename(optional)>\n",
         argv[0]);
179          return -1;
180        }
```

```
181      printf("Command line argument is present...\nLoading configuration data
       ...\n");
182      if (read_macs_and_ports_from_file(argv[1])<0) {
183        printf("File cannnot be opened...\n");
184        return -1;
185      }
186    }
187
188    printf("Create and configure l3 test controller...\n");
189    c = create_controller_with_init(11111, 3, dhf, init);
190    fill_ipv4_fib_lpm_table(ip, port, mac);
191    fill_ipv4_fib_lpm_table(ip2, port2, mac2);
192
193    fill_sendout_table(port, smac);
194    fill_sendout_table(port2, smac);
195
196    printf("Launching controller's main loop...\n");
197    execute_controller(c);
198
199    printf("Destroy controller\n");
200    destroy_controller(c);
201
202    return 0;
203 }
```

Listing C.1 – MACSAD IPv4 Controller code

## C.2   MACSAD IPv6 Controller code

```
1  #include "controller.h"
2  #include "messages.h"
3  #include <unistd.h>
4  #include <stdio.h>
5  #include <string.h>
6  #include <time.h>
7
8  #define MAX_MACS 2000000
9
10 controller c;
11
12 void fill_ipv6_fib_lpm_table(uint8_t ip[16], uint8_t port, uint8_t mac[6])
13 {
14    char buffer[2048];
15    struct p4_header* h;
16    struct p4_add_table_entry* te;
17    struct p4_action* a;
18    struct p4_action_parameter* ap,* ap2;
19    struct p4_field_match_exact* exact;
```

```
20
21    h = create_p4_header(buffer, 0, 2048);
22    te = create_p4_add_table_entry(buffer,0,2048);
23    strcpy(te->table_name, "ipv6_fib_lpm");
24
25    exact = add_p4_field_match_exact(te, 2048);
26    strcpy(exact->header.name, "ipv6.dstAddr");
27    memcpy(exact->bitmap, ip, 16);
28    exact->length = 16*8+0;
29
30    a = add_p4_action(h, 2048);
31    strcpy(a->description.name, "fib_hit_nexthop");
32
33    ap = add_p4_action_parameter(h, a, 2048);
34    strcpy(ap->name, "dmac");
35    memcpy(ap->bitmap, mac, 6);
36    ap->length = 6*8+0;
37
38    ap2 = add_p4_action_parameter(h, a, 2048);
39    strcpy(ap2->name, "port");
40    ap2->bitmap[0] = port;
41    ap2->bitmap[1] = 0;
42    ap2->length = 2*8+0;
43
44    netconv_p4_header(h);
45    netconv_p4_add_table_entry(te);
46    netconv_p4_field_match_exact(exact);
47    netconv_p4_action(a);
48    netconv_p4_action_parameter(ap);
49    netconv_p4_action_parameter(ap2);
50
51    send_p4_msg(c, buffer, 2048);
52 }
53
54 void fill_sendout_table(uint8_t port, uint8_t smac[6])
55 {
56     char buffer[2048];
57     struct p4_header* h;
58     struct p4_add_table_entry* te;
59     struct p4_action* a;
60    struct p4_action_parameter* ap;
61    struct p4_field_match_exact* exact;
62
63     h = create_p4_header(buffer, 0, 2048);
64     te = create_p4_add_table_entry(buffer,0,2048);
65     strcpy(te->table_name, "sendout");
66
```

```c
67          exact = add_p4_field_match_exact(te, 2048);
68          strcpy(exact->header.name, "standard_metadata.egress_port");
69          exact->bitmap[0] = port;
70      exact->bitmap[1] = 0;
71          exact->length = 2*8+0;
72
73      a = add_p4_action(h, 2048);
74      strcpy(a->description.name, "rewrite_src_mac");
75
76      ap = add_p4_action_parameter(h, a, 2048);
77      strcpy(ap->name, "smac");
78      memcpy(ap->bitmap, smac, 6);
79      ap->length = 6*8+0;
80
81          netconv_p4_header(h);
82          netconv_p4_add_table_entry(te);
83          netconv_p4_field_match_exact(exact);
84          netconv_p4_action(a);
85          netconv_p4_action_parameter(ap);
86          send_p4_msg(c, buffer, 2048);
87  }
88
89  uint8_t macs[MAX_MACS][6];
90  uint8_t portmap[MAX_MACS];
91  uint8_t ips[MAX_MACS][16];
92  int mac_count = -1;
93
94  int read_macs_and_ports_from_file(char *filename) {
95      FILE *f;
96      char line[200];
97      int values[6];
98      int values_ip[16];
99      int port;
100     int i;
101
102     f = fopen(filename, "r");
103     if (f == NULL) return -1;
104
105     while (fgets(line, sizeof(line), f)) {
106         line[strlen(line)-1] = '\0';
107         //TODO why %c?
108         if (23 == sscanf(line, "%02x%02x:%02x%02x:%02x%02x:%02x%02x:%02x%02x
        :%02x%02x:%02x%02x:%02x%02x %x:%x:%x:%x:%x:%x %d",
109             &values_ip[0], &values_ip[1], &values_ip[2], &values_ip[3],
110             &values_ip[4], &values_ip[5], &values_ip[6], &values_ip[7],
111             &values_ip[8], &values_ip[9], &values_ip[10], &values_ip[11],
112             &values_ip[12], &values_ip[13], &values_ip[14], &values_ip[15],
```

```
113              &values [0] , &values [1] , &values [2] ,
114              &values [3] , &values [4] , &values [5] , &port) )
115       {
116          if (mac_count==MAX_MACS−1)
117          {
118            printf("Too many entries ...\n");
119            break;
120          }
121          ++mac_count;
122          for( i = 0; i < 6; ++i )
123            macs[mac_count][i] = (uint8_t) values[i];
124          for( i = 0; i < 16; ++i )
125            ips[mac_count][i] = (uint8_t) values_ip[i];
126          portmap[mac_count] = (uint8_t) port;
127        } else {
128          printf("Wrong format error in line %d : %s\n", mac_count+2, line);
129          fclose(f);
130          return −1;
131        }
132      }
133      fclose(f);
134      return 0;
135 }
136
137 void dhf(void* b) {
138      printf("Unknown digest received\n");
139 }
140
141 void init() {
142      int i;
143      uint8_t smac[6] = {0xd0, 0x69, 0x0f, 0xa8, 0x39, 0x90};
144      printf("INIT");
145      clock_t begin = clock();
146      for (i=0;i<=mac_count;++i)
147      {
148        fill_ipv6_fib_lpm_table(ips[i], portmap[i], macs[i]);
149
150              if(0 == (i%1000)){ printf("%d inside sleep \n", i);sleep(1);;}
151        fill_sendout_table(portmap[i], smac);
152              usleep(1000);
153      }
154
155      clock_t end = clock();
156      double time_spent = (double)(end − begin) / CLOCKS_PER_SEC;
157
158          printf ("ctrl Total entries sent %d time %f\n",i,time_spent);
159 }
```

```
160
161    int main(int argc, char* argv[])
162    {
163      uint8_t ip[16] = {0x20,0x01,0x0d,0xb8,0x85,0xa3,0x08,0xd3,0x13,0x19,0x8a
              ,0x2e,0x03,0x70,0x73,0x34};
164        uint8_t mac[6] = {0xa0, 0x36, 0x9f, 0x3e, 0x94, 0xea};
165      uint8_t port = 1;
166
167      uint8_t ip2[16] = {0x13,0x19,0x8a,0x2e,0x03,0x70,0x73,0x35,0x20,0x01,0x0d
              ,0xb8,0xff,0xff,0x08,0xd3};
168        uint8_t mac2[6] = {0xa0, 0x36, 0x9f, 0x3e, 0x94, 0xe8};
169      uint8_t port2 = 0;
170
171        uint8_t smac[6] = {0xd0, 0x69, 0x0f, 0xa8, 0x39, 0x90};
172
173      if (argc >1) {
174        if (argc!=2) {
175          printf("Too many arguments...\nUsage: %s <filename(optional)>\n",
          argv[0]);
176          return -1;
177        }
178        printf("Command line argument is present...\nLoading configuration data
          ...\n");
179        if (read_macs_and_ports_from_file(argv[1])<0) {
180          printf("File cannnot be opened...\n");
181          return -1;
182        }
183      }
184
185      printf("Create and configure l3 test controller...\n");
186      c = create_controller_with_init(11111, 3, dhf, init);
187
188      printf("Launching controller's main loop...\n");
189      execute_controller(c);
190
191      printf("Destroy controller\n");
192      destroy_controller(c);
193
194      return 0;
195    }
```

Listing C.2 – MACSAD IPv6 Controller code

# ANNEX  D  —  P4 Code

## D.1  IPv4 Code

```
1   header_type ethernet_t {
2     fields {
3       dstAddr : 48;
4       srcAddr : 48;
5       etherType : 16;
6     }
7   }
8
9   header_type ipv4_t {
10    fields {
11      versionIhl : 8;
12      diffserv : 8;
13      totalLen : 16;
14      identification : 16;
15      fragOffset : 16;
16      ttl : 8;
17      protocol : 8;
18      hdrChecksum : 16;
19      srcAddr : 32;
20      dstAddr: 32;
21    }
22  }
23
24  parser start {
25    return parse_ethernet;
26  }
27
28  #define ETHERTYPE_IPV4 0x0800
29
30  header ethernet_t ethernet;
31
32  parser parse_ethernet {
33    extract(ethernet);
34    return select(latest.etherType) {
35      ETHERTYPE_IPV4 : parse_ipv4;
36      default: ingress;
37    }
38  }
39
40  header ipv4_t ipv4;
```

```
41
42   parser parse_ipv4 {
43     extract(ipv4);
44     return ingress;
45   }
46
47   action on_miss() {
48   }
49
50   action fib_hit_nexthop(dmac, port) {
51     modify_field(ethernet.dstAddr, dmac);
52     modify_field(standard_metadata.egress_port, port);
53     add_to_field(ipv4.ttl, -1);
54   }
55
56   table ipv4_fib_lpm {
57     reads {
58       ipv4.dstAddr : lpm;
59     }
60     actions {
61       fib_hit_nexthop;
62       on_miss;
63     }
64     size : 512;
65   }
66
67   action rewrite_src_mac(smac) {
68     modify_field(ethernet.srcAddr, smac);
69   }
70
71   table sendout {
72     reads {
73       standard_metadata.egress_port : exact;
74     }
75     actions {
76       on_miss;
77       rewrite_src_mac;
78     }
79     size : 512;
80   }
81
82   control ingress {
83   /* fib lookup, set dst mac and standard_metadata.egress_port */
84     apply(ipv4_fib_lpm);
85
86   /* set smac from standard_metadata.egress_port */
87     apply(sendout);
```

```
88  }
89
90  control egress {
91  }
```

Listing D.1 – IPv4 P4 code

## D.2   IPv6 Code

```
1   header_type ethernet_t {
2     fields {
3       dstAddr : 48;
4       srcAddr : 48;
5       etherType : 16;
6     }
7   }
8
9   header_type ipv6_t {
10      fields {
11          version : 4;
12          trafficClass : 8;
13          flowLabel : 20;
14          payloadLen : 16;
15          nextHdr : 8;
16          hopLimit : 8;
17          srcAddr : 128;
18          dstAddr : 128;
19      }
20  }
21
22  parser start {
23    return parse_ethernet;
24  }
25
26  #define ETHERTYPE_IPV6 0x86DD
27
28  header ethernet_t ethernet;
29
30  parser parse_ethernet {
31    extract(ethernet);
32    return select(latest.etherType) {
33      ETHERTYPE_IPV6 : parse_ipv6;
34      default: ingress;
35    }
36  }
37
38  header ipv6_t ipv6;
39
```

```
40    parser parse_ipv6 {
41       extract(ipv6);
42       return ingress;
43    }
44
45    action on_miss() {
46    }
47
48    action fib_hit_nexthop(dmac, port) {
49       modify_field(ethernet.dstAddr, dmac);
50       modify_field(standard_metadata.egress_port, port);
51       add_to_field(ipv6.hopLimit, -1);
52    }
53
54    table ipv6_fib_lpm {
55       reads {
56          ipv6.dstAddr : lpm;
57       }
58       actions {
59          fib_hit_nexthop;
60          on_miss;
61       }
62       size : 512;
63    }
64
65    action rewrite_src_mac(smac) {
66       modify_field(ethernet.srcAddr, smac);
67    }
68
69    table sendout {
70       reads {
71          standard_metadata.egress_port : exact;
72       }
73       actions {
74          on_miss;
75          rewrite_src_mac;
76       }
77       size : 512;
78    }
79
80    control ingress {
81    /* fib lookup, set dst mac and standard_metadata.egress_port */
82       apply(ipv6_fib_lpm);
83
84    /* set smac from standard_metadata.egress_port */
85       apply(sendout);
86    }
```

```
87
88  control egress {
89  }
```

Listing D.2 – IPv6 P4 code