

OFLOPS-SUME and the art of switch characterization

Rémi Oudin* Gianni Antichi[†], Charalampos Rotsos[‡], Andrew W. Moore[†], Steve Uhlig[§]

*Université Paris-Saclay, [†]University of Cambridge, [‡]Lancaster University, [§]Queen Mary, University of London

Abstract—The philosophy of SDN has introduced new challenges in network device management. In contrast to traditional network devices that contained both the control and the data plane functionality in a tightly coupled manner, SDN technologies separate the two network planes, and define a remote API for low-level device configuration. Nonetheless, the enhanced flexibility of the SDN paradigm is prone to create novel performance and scalability bottlenecks in the network. Especially, the data and control plane performance characteristics of an SDN switch can significantly affect the overall network performance.

To help network managers and application developers better understand the actual behavior of SDN implementations, we present and hardware/software co-design that enables switch characterization at 40Gbps and beyond. We conduct an evaluation of both software and hardware switches. We illustrate the importance of measurements-based SDN switch characterization, by exposing the unwanted effects of the OpenFlow barrier primitive, as well as how simple operations, such as packet modification, can impact the experienced forwarding latency. We release the code publicly as open source to promote experiments reproducibility, as well as encourage the network community to evolve our solution.

Index Terms—OpenFlow, testing, switch performances

I. INTRODUCTION

Research on Software-Defined Networking (SDN) technologies has produced a wide range of applications improving network functionality [1]. As a result, many network vendors, within only a few years, enabled SDN support in their products, in an effort to transfer research innovation into the market [2], [3].

However, SDN adoption has highlighted a series of new network management challenges. Legacy network devices contain both the control and the data plane functionality in a tightly coupled manner. In contrast, SDN technologies define a remote API for low-level device configuration, that clearly separates the control from the data plane. This low-level API is used by Network Operating Systems (NOS) to synthesize high-level network control APIs, which can be used by control applications to realize new network services and applications. The introduction of multiple abstraction layers in SDN network architectures introduces novel performance-critical functional blocks in the control stack, like the control application, the NOS, the switch SDN driver, *i.e.* the interface between the control and data planes in an SDN switch, and the switch silicon.

While much research has focused on control application and NOS benchmarking [4], [5], [6], [7], [8], we argue that the effective deployment of SDN in production networks

requires a flexible and high-precision open-source switch performance characterization platform. The OpenFlow protocol, the predominant SDN realization, currently in its 1.6 version, has greatly morphed since its original versions. This evolution reflects the development and deployment experience of the technology, on a constantly widening range of network environments. Providing an open and flexible solution for SDN testing is important to ascertain that the behavior of specific implementations is in-line with the expectations of the supported OpenFlow versions. The SDN community requires a performance testing platform, capable to co-evolve with the protocol and to support rapid prototyping of experimentation scenarios that test and troubleshoot the impact of new protocol features. Furthermore, as network link speeds continuously increase, meaningful packet-level measurements have higher precision requirements, creating a challenge that can only be met with specialized hardware.

This paper presents the design of OFLOPS-SUME, a NetFPGA-SUME [9] based hardware-software co-design for switch characterization. OFLOPS-SUME builds upon OFLOPS [10], a software SDN testing platform, and OSNT [11], a high-performance traffic generation and measurement platform for the NetFPGA SUME. We introduce the new design of OSNT, previously available only on the NetFPGA-10G board [12], as well as its integration with OFLOPS. We use the proposed architecture to test and characterize both software and hardware OpenFlow switches. To demonstrate the relevance of OFLOPS-SUME, we use it to expose the undesirable impact of the barrier primitive, as well as the impact of simple packet manipulation operations on the observed switch performance. OFLOPS-SUME can be used with any hardware/software device offering an OpenFlow control plane, including P4-programmable [13] data planes with OpenFlow support. Moreover, by open-sourcing the tool, we invite the community to develop support for additional SDN control planes, such as the P4-runtime [14].

The contributions of the paper can be summarized as follows:

- We describe the architecture of OSNT on NetFPGA-SUME as well as its integration with OFLOPS, to obtain a high performance, 40Gbps and beyond, OpenFlow switch testing and characterization platform.
- Using OFLOPS-SUME, we conduct a performance study of both software and hardware switches and discuss the implications of the different switch profiles on the overall network behavior, policy consistency, scalability

and performance.

- We release the code publicly as open source¹. By providing an open source solution, we promote experiments reproducibility as well as encourage the network community to adapt, evolve, and improve our solution to specific needs.

The remainder of this paper is organized as follows: we first motivate the need for our switch characterization tool by discussing challenges and opportunities in network testing (§II). We then describe the architecture of OFLOPS-SUME (§III) and provide a measurement study of both software and hardware switches (§IV). Finally, we present related efforts (§V) and conclude our work (§VI).

II. CHALLENGES AND OPPORTUNITIES IN NETWORK TESTING

Two major drivers for SDN are open APIs and the network management flexibility. These two aspects have incentivised the community to devise protocols and standards for the control of multi-technology and multi-vendor networks [15]. However, implementing them is challenging in practice, due to the wide diversity of the actual support and performance across different vendors.

A good example of how an SDN implementation can impact network operations can be found in the blackholing practice: a defense strategy against Distributed Denial of Service (DDoS) attacks, which blocks malicious traffic by redirecting it into a black hole or null route. Characterizing the network control responsiveness, by quantifying *the latency to update the forwarding state of k flows*, is crucial to understand when exactly a blackhole policy is active. Furthermore, the performance of a SDN implementation can introduce transients policy violations during the deployment of a network policy, due to slow or re-ordered rule updates [16]. To ensure consistent policy enforcement, individual devices must implement a specific rule set at all times, even during policy updates. Rule update ordering in OpenFlow networks can be guaranteed using Barrier messages. However, disambiguation in Barrier message semantics between commercial switches can result in unexpected behaviors [17]. Understanding *how accurately a barrier message represents the data plane configuration status* can thus help in implementing better update strategies that overcome policy inconsistencies. Finally, the deployment of hybrid Xeon and FPGA servers in datacenters [18] enables a wide range of performance improvements, allowing a system to control the trade-off between performance (FPGA) and flexibility/scalability (CPU).

A. Switch Architecture

This section provides a high-level design overview of a generic *switch*. We highlight the physical limits in control and data plane performance and motivate our discussion on performance characterization. Our presentation focuses on top-of-rack (ToR) switch devices and high-performance software

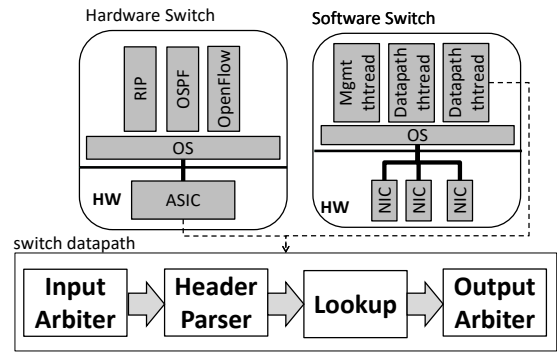


Fig. 1: Functional blocks of a network switch.

switch frameworks, the most common network device types with built-in SDN support. The key functional components are depicted in Figure 1.

Switch datapath. The datapath usually consists of four functional blocks. The *input arbiter* module schedules incoming packets to the packet processing threads, for software switches, or to the main processing pipeline, for hardware switches. The main processing pipeline is implemented by the *header parser* and the *lookup* modules. The parser module consists of multiple protocol parsers that extract important header information, used by the lookup module to define the per-packet processing and forwarding policy. The header parsing and lookup modules can recirculate packets through the processing pipeline, to parse nested packet headers or to enable multi-table processing datapaths. Finally, the *output arbiter* module enforces traffic policing, applies outstanding packet modifications and forwards packets to the appropriate output ports.

Modern hardware based ToR switches use either an ASIC [19], [20], a FPGA [21] or a Network Processor unit [22] for packet processing acceleration. In all cases, the main bottleneck is the size and access speed of its lookup memory modules. In contrast, software switches rely on the parallelization of modern multi-core CPUs to achieve line-rate performance. In this scenario, the main bottleneck is the number of cores and the clock rate. Furthermore, software switches use PCIe channels to connect the system CPUs with the server NICs. The speed and capacity of the PCIe interconnect and the card driver performance can thus be a bottleneck.

Switch Management. The switch management module is responsible for translating high-level control operations into appropriate data plane configuration. In hardware switches, the switch management module is made of different control agents, each providing support for a specific control function, such as legacy routing protocol or SDN support. Depending on the switch model, to support the computational requirements of the control agents, network vendors typically use a general-purpose low-power management CPU, such as PowerPC SoC and Intel Atom. In the former case, the CPU may become a bottleneck when used by interactive control applications. In

¹<https://github.com/oflops-nf>

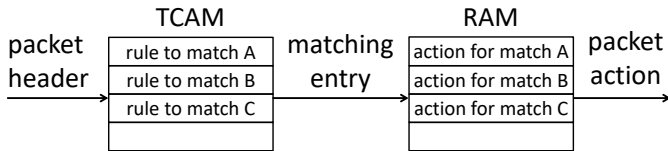


Fig. 2: The two memories used by hardware switch datapaths: (i) TCAMs for flow-match definitions, and (ii) RAMs for forwarding actions.

contrast, the switch management module in software switch runs a separate thread or process and uses OS inter-process communication mechanisms to configure the state of the packet processing threads. Because the management module runs on the same CPU as the packet processing threads, software switch control channels offer enhanced responsiveness.

Switch Datapath-Management interconnect. The management modules typically use a dedicated PCIe interconnect to control and monitor the hardware unit. The latency and speed of this interconnect define the responsiveness of the control channel in hardware switches. The newest Broadcom Trident 3 chipset [23] uses a PCIe Gen3 x4 channel, providing a theoretical maximum capacity of approximately 4Gbps.

B. Datapath memory configuration in modern SDN switches

The packet processing functionality in modern SDN switches relies on a fast-lookup abstraction with support for extensive wildcard matching and reconfigurability. Such an abstraction is based on a match-action architecture [24].

Software switches typically implement the match-action tables using software lookup structures, like tuple space search classifiers [25]. Updating such structures usually has very low latency, but the lookup speed depends on the size and the occupancy of the lookup structure. In contrast, hardware switches implement tables using two different memory subsystems embedded in the hardware. The match definitions are stored in one or more Ternary Content-Addressable Memory (TCAM) modules. TCAMs are optimized for fast lookups, supporting $O(1)$ lookup complexity for any match, irrespective of the lookup key-width or the number of stored match definitions. The output of a TCAM lookup is then used to index a RAM memory module, which stores the actions and statistics associated to the match [26], [27], [24]. Figure 2 shows the interaction between TCAMs and RAM memory in high performance switching engines.

The management module changes the per-packet processing behavior of the switch by manipulating the two memories. Adding a new rule causes an update operation in both TCAM and RAM memories. A modification of the behavior of an existing rule causes an update operation in the RAM only. Unfortunately, the update latency of a single TCAM entry is not constant and depends on the memory design and on the number and structure of the entries already installed [16]. The update latency of a RAM entry is constant and depends solely on the memory type. This update latency mismatch can cause inconsistencies during update operations.

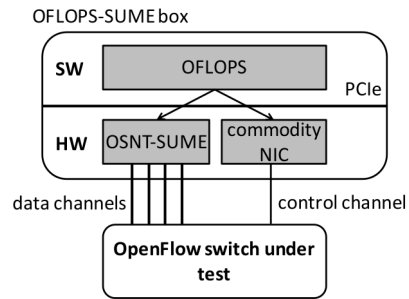


Fig. 3: OFLOPS-SUME design.

III. OFLOPS-SUME: ARCHITECTURE

A schematic of the OFLOPS-SUME design is depicted in Figure 3. The OFLOPS-SUME architecture consists of a software and a hardware subsystem. The software subsystem runs the core OFLOPS functionality, along with the test of the user. A test contains both the control and data plane logic of the experiment.

While the control logic, i.e., OpenFlow control messages, interacts with the system under test through a commodity network interface card, the data plane logic relies on a NetFPGA-SUME (running the OSNT design) hardware subsystem to fulfill the data plane requirements of the experiment.

We chose to use two separate hardware subsystems for control and data plane channels to balance flexibility and performance. While control plane messages, i.e., OF messages, mainly works at milli-seconds timescales and do not require hardware acceleration, data plane traffic needs high-throughput generation and hardware timestamping to properly characterize switches at 10Gbps and beyond. Control messages are therefore handled by a commodity network interface card, so that future updates of the control message protocol require a software update only. This makes the design flexible and future-proof. On the other hand, data plane traffic is handled using the OSNT hardware acceleration, to guarantee performance.

The integration of OFLOPS and OSNT on NetFPGA-SUME is achieved through a C library², which exposes a traffic generation and capture interface, as well as access to card statistics (e.g., counter for packet captures and drops). The library is designed to exploit as much as possible the OSNT capabilities, i.e., multi-10Gbps packet generation and nanosecond scale hardware timestamping. The user can interconnect the OFLOPS-SUME host with one or more switches in arbitrary topologies, and measure with high precision specific aspects of the network architecture, both of the data and control plane. Next, we describe the design of the OSNT on NetFPGA-SUME hardware (Section III-A) and the OFLOPS software architecture (Section III-B).

A. OSNT: Open Source Network Tester

OSNT is an open-source hardware-based traffic generation and capture system. Originally designed for the research and

²<https://github.com/oflops-nf/nf-pktgencap-lib>

teaching community, its key design goals were low cost, high-precision time-stamping and packet transmission, as well as scalability. Being open-source, it provides the flexibility to allow new protocol tests to be added. The first prototype implementation was built upon the NetFPGA-10G platform, while the currently supported version relies on the newer NetFPGA-SUME hardware.

Figure 4 shows the current OSNT architecture. While the hardware side of OSNT (*i.e.*, NetFPGA-SUME board) provides two independent pipelines for high-precision traffic generation and monitoring, the software provides user APIs to interact with the hardware through the NetFPGA driver.

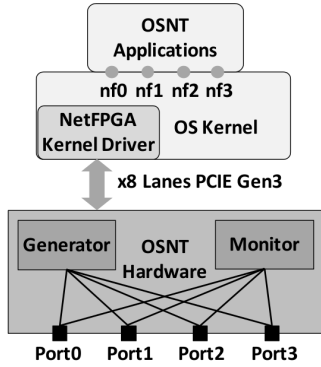


Fig. 4: Architecture of OSNT.

In the following, we first describe the timestamping operation (Section III-A1), to proceed with the architectures of the traffic monitor (Section III-A2) and traffic generation (Section III-A3) pipelines.

1) *Timestamping*: Providing an accurate timestamp to packets is a fundamental objective of OSNT. Packets are timestamped as close to the physical Ethernet device as possible, so as to minimize FIFO-generated jitter and obtain accurate latency measurements. A dedicated timestamping unit, shared between the monitoring and the generation pipelines, stamps packets as they arrive from/to the physical (MAC) interfaces.

To minimize overhead while also providing sufficient resolution and long-term stability, OSNT uses a 64-bit timestamp divided into two parts: the upper 32-bits count seconds, while the lower 32-bits provide a fraction of a second with a maximum resolution of about 233ps. Given the 160Mhz board clock, the practical resolution is 6.25ns. Note that accurate timekeeping requires correcting the frequency drift of an oscillator. To this end, OSNT uses Direct Digital Synthesis (DDS), a technique by which arbitrary variable frequencies can be generated using synchronous digital logic [28]. The addition of a stable pulse-per-second (PPS) signal, *e.g.*, derived from a GPS receiver, permits both high long-term accuracy and the synchronization of multiple OSNT elements.

2) *Traffic Monitor*: The OSNT traffic monitor subsystem provides functions for packet capturing, hardware packet filtering, high precision packet timestamping and high-level traffic statistic gathering.

Figure 5 illustrates the architecture of the monitoring pipeline. A module placed immediately after the physical interfaces, and before the receive queues, timestamps incoming packets as they are received by the hardware. The four per-port pipelines are then merged into a single one using an arbitration process. While packets are temporarily stored in a small packet FIFO (left pipeline after the arbiter), awaiting for a decision, their copy is also sent to the right pipeline. As soon as they hit the *statistic collector*, internal hardware counters exposed to the software are updated accordingly to provide high-level statistics about the traffic. Then, the *header extractor* module gathers the relevant information (the 5-tuple³) to feed the *TCAM*. Indeed, even though our design implicitly copes with a workload of full line-rate per port of minimum sized packets, the input traffic will often exceed the capacity of the host-processing, storage, etc., or may not be of practical interest. To this end, OSNT implements two traffic-thinning approaches:

- 5-tuple filter implemented with a TCAM;
- Packet cutting.

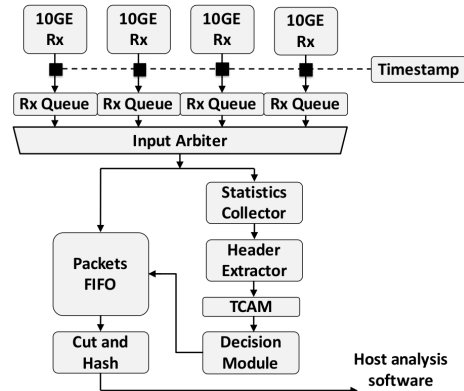


Fig. 5: Architecture of the OSNT traffic monitoring system.

The TCAM, implemented using FPGA resources, enables the host to write up to sixteen 5-tuple rules representing the traffic of interest. Only packets that are matched in the TCAM are sent to the software, while all other packets are discarded. Once the TCAM provides the matching result, the decision module informs the FIFO about the action to be done on the packet: drop or forward. If the packet is not dropped, it will pass through the *cut and hash* module before hitting the host. This last module, configurable from the software, allows to pass only a fixed-length part of each packet (sometimes called as *snap-length*) along with a hash of the entire original packet.

3) *Traffic Generation*: The OSNT traffic generator pipeline is designed to generate 10Gbps at line-rate, 64B packets, per card interface. Currently, it is able to replay network traffic from PCAP files stored in the NetFPGA internal or external memories.

Figure 6 illustrates its high-level architecture. The *Generator management* is a set of software APIs that allows the host

³We define 5-tuple as the combination of IP address pair, layer four port pair and protocol.

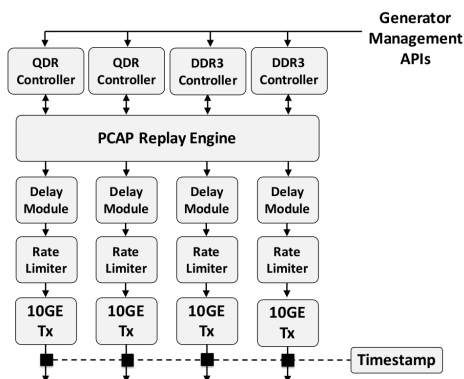


Fig. 6: Architecture of the OSNT traffic generation system.

to interact with the hardware, to configure the generation parameters. With such APIs, the host can record in the NetFPGA external memories, i.e., SRAM and DRAM, the PCAP traces to be reproduced. Each bank of memory is associated to a specific physical port, to avoid cross dependency between the generated traffic. However, such a solution limits the amount of data that can be stored (and then generated) for each port. Specifically, the 2 ports coupled with the SRAM can store a maximum of 14MB of data each, while the ones assigned to the DRAM can store a maximum of 4GB each. The *PCAP replay Engine*, upon signaling from the software, reads the external memory and passes the packet along the pipeline. The generation process can follow user-defined timing information, i.e., inter-packet gap, or can strictly depend on the PCAP trace. In the former case, the packets are delivered to a *Delay Module* which assigns a configured fixed delay for each packet and then to a *Rate Limiter* to guarantee that the generation does not exceed the assigned rate. In the latter case, the Delay Module and the Rate Limiter are bypassed, as the PCAP trace has already enough timing information to proceed with the generation. Finally, the traffic generator has an accurate timestamping mechanism, located just before the transmit 10GbE MAC. The mechanism, identical to the one used in the traffic monitoring unit, is used for timing-related measurements of the network, allowing characterization of measurements such as latency and jitter. The timestamp is embedded within the packet alongside a packet counter at a pre-configured location, and can be extracted at the receiver as required.

B. OFLOPS: OpenFlow Operations Per Second

OFLOPS is a modular software framework enabling precise switch performance characterization. It provides an environment for network experimenters to develop interaction scenarios with OpenFlow switches and monitor the switch behavior over multiple control and information channels. The platform offers a unified API to control the data, the control and the management plane logic, as well as log low-level network measurement data. Experimenters can use the API to realize the control and monitoring logic of an interaction scenario, and describe the flow-level traffic characteristics. The

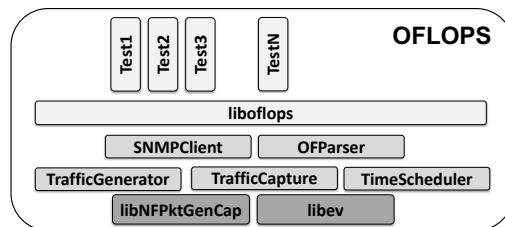


Fig. 7: OFLOPS platform architecture

code of an experiment is compiled into a dynamic libraries, which will be loaded at run-time by the OFLOPS-runtime based on a configuration file. The OFLOPS traffic monitor and generation API abstracts the complexity to realize the respective functionalities across different platforms.

OFLOPS follows a 3-layer architecture, depicted in Figure 7. The bottom layer consists of a set of services enabling interaction with a switch across different control and data channels. Access to the different information and control channels is unified and abstracted by the middle layer of the platform. Finally, the top layer of the OFLOPS architecture consists of measurement modules which implement the interaction scenarios.

In order to precisely monitor the data plane of modern 10 GbE switches, OFLOPS must be able to process multi-million packets per second and offer sub- μ sec measurement precision. As a result, the OFLOPS platform is designed to minimize packet processing overheads and uses a multi-threaded architecture to parallelize event processing, in an effort to minimize packet loss and measurement noise. At boot time, the process initializes five threads, each responsible of managing specific event types, namely (i) traffic generation, (ii) traffic capture, (iii) OF message parsing, (iv) SNMP responses and (v) time events. To avoid synchronization overheads, OFLOPS provides a *non* thread-safe event API. Module developers are responsible of implementing custom synchronization mechanisms within their test logic.

OFLOPS traffic generation and capture supports a wide range of back-end mechanisms, including the OSNT-SUME hardware design, and the pktgen and libpcap kernel modules. The different back-ends allow experimenters to test switches using commodity hardware, while maximizing measurement precision using specialized hardware equipment. The same testing code can use programmable hardware to offload high volume packet processing complexity in a hardware card, or realize traffic generation in software in resource-limited measurement setups. Currently, OFLOPS supports the most widely adopted OF versions: 1.0 and 1.3. The OFLOPS code base contains a set of reference test modules supporting all elementary protocol interactions, like flow table manipulation and traffic monitoring.

IV. PERFORMANCE CHARACTERIZATION

In this section, we present a switch characterization study, using the OFLOPS-SUME platform. In our analysis, we employ three representative OpenFlow switch types based on the

OpenVSwitch (OVS) [29] switch platform. Firstly, we employ the software kernel-based OVS (OVS-kernel) switch, highly popular in cloud infrastructures. Secondly, we employ a port of the OVS platform on the DPDK framework (OVS-DPDK)⁴, a high-performance zero-copy packet processing framework used extensively to implement network function applications. Finally, we used a white-box hardware switch (EdgeCore AS5812-54X) running an OpenFlow-enabled firmware (PicaOS), a common design choice for ToR switches. The switch offers an Intel Atom quad-core co-processor, which runs a modified OVS version capable to offload data plane processing to a Broadcom ASIC (BCM56854) supporting up to 1.28Tbps throughput and a 32k entry TCAM module.

Our measurement study focuses on switch forwarding performance (§ IV-A), the flow table management capabilities, and the consistency characteristics during forwarding policy reconfiguration (§ IV-B). Figure 3 depicts our measurement setup. The OFLOPS-SUME runs on an Ubuntu server (quad-core Intel E5-1603, 16 GB RAM) equipped with a NetFPGA-SUME card and programmed with the OSNT bitfile. The software OVS instances (v2.8.1) run on an Ubuntu server (16-core Intel E5-2630, 32 GB RAM) equipped with a dual-port 10G Intel NIC (82599ES). For all the figures, unless stated otherwise, we present the median value over 20 experiment repetitions alongside the minimum and the 90th percentile.

A. Data plane performance

Figure 8 depicts the forwarding latency of each OpenFlow switch for varying packet sizes. Figure 8a reports the results for the Edgecore, when no specific packet modification is requested, *i.e.*, only simple forwarding. We observe that packet size variations have a negligible impact on the median latency, which remains stable around 2.5 μ s. Nonetheless, increasing packet sizes inflate significantly the 90th percentile of the forwarding latency, which is 66% higher for maximum size packets in comparison to 128-byte sized packet.

Software switch forwarding latency, depicted in Figure 8b, is approximately one order of magnitude higher than the hardware switch. In addition, OVS-Kernel exhibits higher latency variability than OVS-DPDK. These latency differences can be attributed to the underlying packet processing architecture and, specifically, the mechanism transferring packets between the CPU and the NIC, as demonstrated in Figure 8c. The figure depicts the increase of the median latency caused by a packet manipulation operation⁵. While the EdgeCore (not shown in the figure) and OVS-DPDK remain unaffected, OVS-Kernel experiences a median latency increase of 12 μ s in the worst case. Such forwarding latencies have been reported to significantly affect the performance of several common data-center applications [30]. The impact of packet modification is more pronounced for small packets. This can be attributed to the NAPI interrupt mitigation mechanism used by the OVS-Kernel switch, which switches between interrupt-driven and

polling-based modes depending on the incoming packet rate. In contrast, OVS-DPDK uses exclusively a poll-based packet fetching mechanism, which minimizes the impact of the packet size on forwarding latency, but increases the CPU utilization.

B. Control plane performance

Figure 9 shows the results of our control plane performance analysis. We characterize the amount of time needed to insert new flow rules or modify existing ones. In addition, we also focus our attention on assessing the reliability of the barrier primitive, compared to the actual rule installation.

Figure 9a presents the rule update latency of the Edgecore switch for a varying number of rules, using two measurement metrics. The *control plane* latency represents the time between the `FLOW_MOD` message transmission and the reception of the `BARRIER_REPLY` by the switch through the control plane channel (see Figure 3). In contrast, the *data plane* latency measures the time between the `FLOW_MOD` message transmission and the time when at least one packet has been received with the new configuration by OFLOPS-SUME on one of the data channels (see Figure 3). The significant difference between the two measurements (Figure 9a) can be attributed to the behavior of the switch OpenFlow agent. Specifically, we notice that the switch OpenFlow agent transmits the `BARRIER_REPLY` as soon as the flows are received and processed by the co-processor, without checking if the updates have been applied to the ASIC lookup modules. The reported result has a major implication: for a few seconds, control and data plane experience a policy inconsistency which can lead to thousands of mis-routed packets in high performance networks.

Figure 9b shows the performance of the Edgecore switch, when we insert an increasing number of rules. Similarly to the flow modifications, the discrepancy between control and data plane is pretty clear. In addition, the increased rule insertion latency, in comparison to flow modifications, can be attributed to the TCAM reordering operation, required to effectively manage TCAM entries. Since our test inserts rule batches with an increasing priority, the switch OpenFlow agent has to re-order the entries every time we add new TCAM rules, to guarantee correct hardware behavior, thus increasing the switch reconfiguration latency. We do not observe a similar behavior for flow modifications (Figure 9a), since such operations require switch RAM memory writes only and they do not require any TCAM modifications (§II-B). Moreover, it is worth noting the big difference between flow modification and flow addition latencies. The former, for 1400 rules, is approximately 70% faster than the latter. This is again a consequence of the TCAM reordering: a single rule installation can cause many writes in the TCAM and associated RAM.

Finally, Figure 9c shows the behavior of rule addition latency for a software OpenFlow switch. Note the almost three orders of magnitude difference with respect to the hardware case. We omit the flow modification results, since they exhibit the same pattern as the flow additions. Effectively, software switches use in-memory data structures to perform the lookup process which have a low update complexity. This

⁴<http://docs.openvswitch.org/en/latest/intro/install/dpdk/>

⁵In the test, we considered a layer 4 port rewriting.

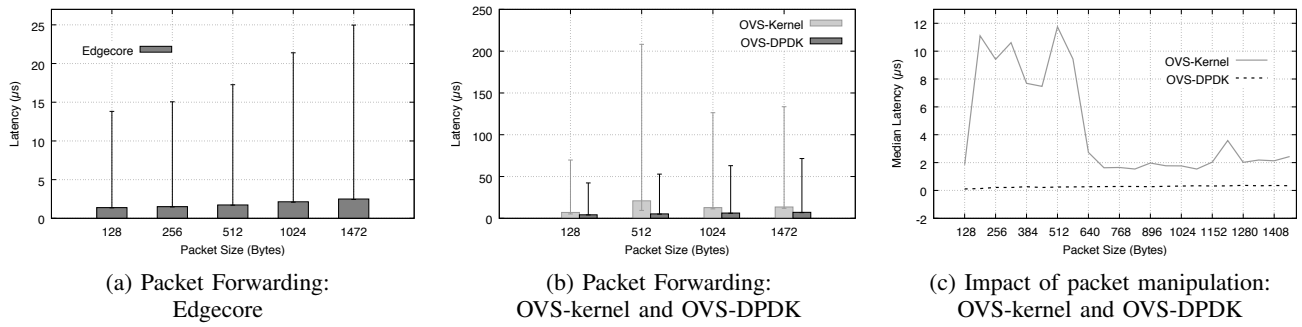


Fig. 8: Data plane performance analysis for Edgecore, OVS-kernel and OVS-DPDK.

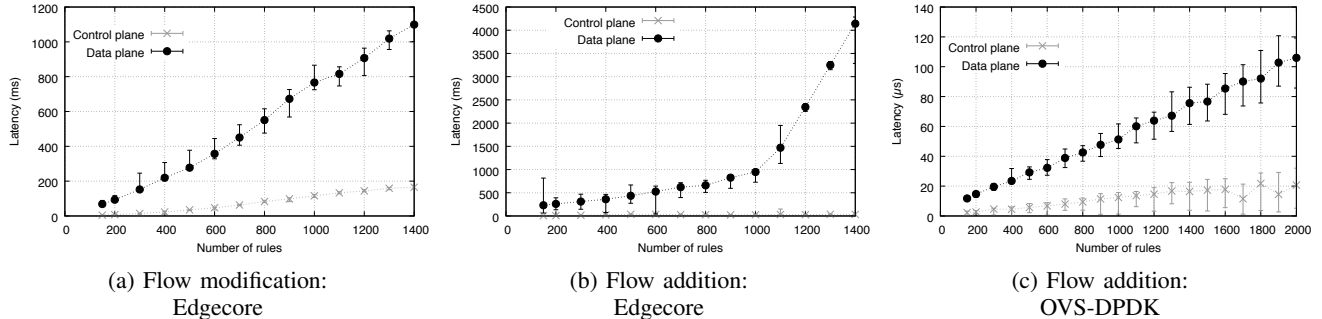


Fig. 9: Control plane performance analysis for Edgecore, OVS-kernel and OVS-DPDK.

makes the rule addition or rule modification processes almost similar: the latency increases linearly with the number of rules changed. In contrast, as discussed in the challenges section (§II), hardware solutions employ TCAMs and RAMs lookup modules which exhibit different update latency characteristics and lead to different behavior during flow additions and flow modifications.

V. RELATED WORK

The SDN community has developed a wide range of validation and testing tools. OFtest [31] is a standards conformity evaluation framework for OpenFlow switches. The suite offers tests to verify the conformance of a switch implementation against a specific OpenFlow standard. In parallel, Cbench [6], OFLOPS [10] and OFLOPS-TURBO [32] are focused primarily on the performance side. Cbench evaluates the performance of a controller by emulating a number of switches generating PKT_IN messages and measuring its responsiveness. The OFLOPS and OFLOPS-Turbo platforms concentrate on switch performance evaluation. While OFLOPS is a pure software solution limited to 1 GbE devices and supporting the 1.0 OpenFlow protocol version, OFLOPS-Turbo integrates an old version of NetFPGA to support higher throughput. However, its limited hardware resources make the solution unattractive to enable the support of new OpenFlow versions as well as more refined data plane tests.

Switch characterization studies have also aimed to identify the impact of SDN technologies on network performance. Bianco *et al.* [33] conducted the first data plane performance

study of software OpenFlow switches. In their study, they highlighted the performance improvements of Open vSwitch over traditional Linux bridging, and showed the performance improvement of NAPI interrupt mitigation on network performance. Starting from their lessons learned, in this paper, we also show the impact of emerging fast I/O frameworks, such as DPDK, on software switching performance. Huang *et al.* [34] benchmarked OpenFlow-enabled switches and illustrated how their implementation can significantly impact data plane latency and throughput. They also presented a measurement methodology and emulator extension, to reproduce these control path performance characteristics, restoring the fidelity of emulation. While the authors focused primarily on vendor-specific variations in the control plane, in this paper, we propose a benchmark for both control and data plane, taking advantage of the OSNT system. Tango [35] is a controller design capable of evaluating the performance profile of a network switch. The controller uses control plane traffic injection and capture during a training phase, to evaluate the performance characteristics of a network switch. The controller uses these results at run-time to adapt the policy deployment strategy. Unfortunately, the precision of the resulting system remains poor, as the control channel offers pretty limited performance and precision guarantees. Furthermore, Curtis *et al.* [36] discussed the switch specific limitations to support OpenFlow functionality in large network deployments. Their measurement study highlighted the impact of the channel between the co-processor and the ASIC in policy reconfiguration. Finally, Jarschel *et al.* [37] employed a DAG-

based network environment to measure the flow manipulation performance of OpenFlow switches, and presented a markov model for the forwarding speed and blocking probability of OpenFlow devices.

VI. CONCLUSIONS

We presented OFLOPS-SUME, a performance characterization platform for OpenFlow switches. OFLOPS-SUME combines the advanced hardware capabilities of OSNT, implemented on the NetFPGA-SUME, with the extensibility of the OFLOPS software framework to enable switch characterization at 40Gbps and beyond.

Using OFLOPS-SUME, we evaluated both software and hardware switches. The performed measurement campaign exposes the sometimes undesirable impact of barrier primitives, as well as the cost in terms of added latency of simple operations, such as packet modification. OFLOPS-SUME can be used with any hardware/software device that uses OpenFlow as a control plane. Thus, it can be also exploited by highly programmable devices, such as P4-enabled ones, as long as they are managed with an OpenFlow control plane. We release the code publicly as open source to promote experiments reproducibility, as well as to encourage the network community to evolve our solution.

REFERENCES

- [1] D. Kreutz, F. M. V. Ramos, P. Veríssimo, C. E. Rothenberg, S. Azodolmoly, and S. Uhlig, "Software-Defined Networking: A Comprehensive Survey," in *Proceedings of IEEE, Volume: 103, Issue: 1*. IEEE, 2015.
- [2] S. Jain, A. Kumar, S. Mandal, J. Ong, L. Poutievski, A. Singh, S. Venkata, J. Wanderer, J. Zhou, M. Zhu, J. Zolla, U. Hölzle, S. Stuart, and A. Vahdat, "B4: Experience with a Globally-deployed Software Defined Wan," in *Special Interest Group on Data Communication (SIGCOMM)*. ACM, 2013.
- [3] M. Kobayashi, S. Seetharaman, G. Parulkar, G. Appenzeller, J. Little, J. Van Reijendam, P. Weissmann, and N. McKeown, "Maturing of OpenFlow and Software-defined Networking Through Deployments," in *Computer Network, Volume: 61*. Elsevier, 2014.
- [4] M. Canini, D. Venzano, P. Perešini, D. Kostić, and J. Rexford, "A NICE Way to Test Openflow Applications," in *Networked Systems Design and Implementation (NSDI)*. USENIX, 2012.
- [5] M. Jarschel, F. Lehrieder, Z. Magyari, and R. Pries, "A Flexible OpenFlow-Controller Benchmark," in *European Workshop on Software Defined Networking (EWSN)*. IEEE, 2012.
- [6] A. Tootoonchian, S. Gorbunov, Y. Ganjali, M. Casado, and R. Sherwood, "On Controller Performance in Software-defined Networks," in *Hot Topics in Management of Internet, Cloud, and Enterprise Networks and Services (Hot-ICE)*. USENIX, 2012.
- [7] D. Turull, M. Hidell, and P. Sjodin, "Performance evaluation of OpenFlow controllers for network virtualization," in *High Performance Switching and Routing (HPSR)*. IEEE, 2014.
- [8] Y. Zhao, L. Iannone, and M. Riguidel, "On the Performance of SDN Controllers: A Reality Check," in *Network Function Virtualization and Software Defined Network (NFV-SDN)*. IEEE, 2015.
- [9] N. Zilberman, Y. Audzevich, A. G. Covington, and A. W. Moore, "NetFPGA SUME: Toward 100 Gbps as Research Commodity," in *Micro, Volume: 34, Issue: 5*. IEEE, 2014.
- [10] C. Rotsos, N. Sarrar, S. Uhlig, R. Sherwood, and A. W. Moore, "OFLOPS: An Open Framework for Openflow Switch Evaluation," in *Passive and Active Measurement (PAM)*. Springer-Verlag, 2012.
- [11] G. Antichi, M. Shahbaz, Y. Geng, N. Zilberman, A. G. Covington, M. Bruyere, N. McKeown, N. Feamster, B. Felderman, M. Blott, A. W. Moore, and P. Owezarski, "OSNT: Open Source Network Tester," in *Network, Volume: 28, Issue: 5*. IEEE, 2014.
- [12] M. Blott, "FPGA Research Design Platform Fuels Network Advances," in *Xcell, Issue: 73*. Xilinx, 2010.
- [13] P. Bosshart, D. Daly, G. Gibb, M. Izzard, N. McKeown, J. Rexford, C. Schlesinger, D. Talayco, A. Vahdat, G. Varghese, and D. Walker, "P4: Programming Protocol-independent Packet Processors," in *SIGCOMM Computer Communication Review, Volume: 44, Issue: 3*. ACM, 2014.
- [14] P. L. Consortium, "P4runtime," <https://p4.org/p4-runtime>.
- [15] M. Casado, N. Foster, and A. Guha, "Abstractions for Software-defined Networks," in *Communications of the ACM, Volume: 57, Issue: 10*. ACM, 2014.
- [16] H. Chen and T. Benson, "The Case for Making Tight Control Plane Latency Guarantees in SDN Switches," in *Symposium on SDN Research (SOSR)*. ACM, 2017.
- [17] J. H. Han, P. Mundkur, C. Rotsos, G. Antichi, N. Dave, A. W. Moore, and P. G. Neumann, "Blueswitch: enabling provably consistent configuration of network switches," in *Architectures for Networking and Communications Systems (ANCS)*. ACM/IEEE, 2013.
- [18] P. K. Gupta, "Xeon+FPGA Platform for the Data Center," <https://www.ece.cmu.edu/~calcm/carl/lib/exe/patch.php?media=carl15-gupta.pdf>, online; accessed August 2017.
- [19] "Juniper high-performance data center switches," <https://www.juniper.net/uk/en/products-services/switching/qfx-series/>, online; accessed February 2018.
- [20] "Cisco data center switches," <https://www.cisco.com/c/en/us/products/switches/data-center-switches/index.html>, online; accessed February 2018.
- [21] "Corsa Technology," www.corsa.com/solutions/, online; accessed February 2018.
- [22] "NoviFlow Top-Of-Rack switching," <https://noviflow.com/tor/>, online; accessed February 2018.
- [23] "Broadcom Trident 3," <https://www.broadcom.com/blog/new-trident-3-switch-delivers-smarter-programmability-for-enterp>, online; accessed February 2018.
- [24] P. Bosshart, G. Gibb, H.-S. Kim, G. Varghese, N. McKeown, M. Izzard, F. Mujica, and M. Horowitz, "Forwarding Metamorphosis: Fast Programmable Match-action Processing in Hardware for SDN," in *Special Interest Group on Data Communication (SIGCOMM)*. ACM, 2013.
- [25] V. Srinivasan and G. Varghese, "Fast Address Lookups Using Controlled Prefix Expansion," in *Transactions on Computer Systems, Volume: 17, Issue: 1*. ACM, 1999.
- [26] R. Giladi, "Network Processors: Architecture, Programming and Implementation." Elsevier, 2008.
- [27] N. Stringfield, R. White, and S. McKee, "Cisco Express Forwarding." Cisco Press, 2013.
- [28] P. Saul, "Direct digital synthesis," in *Circuits and Systems Tutorials*. IEEE Press, 1996.
- [29] B. Pfaff, J. Pettit, T. Koponen, E. J. Jackson, A. Zhou, J. Rajahalme, J. Gross, A. Wang, J. Stringer, P. Shelar, K. Amidon, and M. Casado, "The Design and Implementation of Open vSwitch," in *Networked Systems Design and Implementation (NSDI)*. USENIX, 2015.
- [30] N. Zilberman, M. Grosvenor, D.-A. Popescu, N. Manihatty-Bojan, G. Antichi, M. Wojcik, and A. W. Moore, "Where Has My Time Gone?" in *Passive and Active Measurement (PAM)*. Springer, 2017.
- [31] P. Floodlight, "oftest," <http://www.projectfloodlight.org/oftest/>.
- [32] C. Rotsos, G. Antichi, M. Bruyere, P. Owezarski, and A. W. Moore, "OFLOPS-Turbo: Testing the next generation OpenFlow switches," in *International Conference on Communications (ICC)*. IEEE, 2015.
- [33] A. Bianco, R. Birke, L. Giraud, and M. Palacin, "OpenFlow switching: Data plane performance," in *International Conference on Communications (ICC)*. IEEE, 2010.
- [34] D. Y. Huang, K. Yocum, and A. C. Snoeren, "High-fidelity Switch Models for Software-defined Network Emulation," in *Hot Topics in Software Defined Networking (HotSDN)*. ACM, 2013.
- [35] A. Lazaris, D. Tahara, X. Huang, E. Li, A. Voellmy, Y. R. Yang, and M. Yu, "Tango: Simplifying SDN Control with Automatic Switch Property Inference, Abstraction, and Optimization," in *Conference on Emerging Networking Experiments and Technologies (CoNEXT)*. ACM, 2014.
- [36] A. R. Curtis, J. C. Mogul, J. Tourrilhes, P. Yalagandula, P. Sharma, and S. Banerjee, "DevoFlow: Scaling flow management for high-performance networks," in *Special Interest Group on Data Communication (SIGCOMM)*. ACM, 2011.
- [37] M. Jarschel, S. Oechsner, D. Schlosser, R. Pries, S. Goll, and P. Tran-Gia, "Modeling and Performance Evaluation of an OpenFlow Architecture," in *International Teletraffic Congress (ITC)*. IEEE, 2011.