

PISCES: A Programmable, Protocol-Independent Software Switch

Muhammad Shahbaz*, Sean Choi[◇], Ben Pfaff[†], Changhoon Kim[‡],
Nick Feamster*, Nick McKeown[◇], Jennifer Rexford*

*Princeton University [◇]Stanford University [†]VMware, Inc [‡]Barefoot Networks, Inc
<http://pisc.es.cs.princeton.edu>

Abstract

Hypervisors use software switches to steer packets to and from virtual machines (VMs). These switches frequently need upgrading and customization—to support new protocol headers or encapsulations for tunneling and overlays, to improve measurement and debugging features, and even to add middlebox-like functions. Software switches are typically based on a large body of code, including kernel code, and changing the switch is a formidable undertaking requiring domain mastery of network protocol design *and* developing, testing, and maintaining a large, complex codebase. Changing how a software switch forwards packets should not require intimate knowledge of its implementation. Instead, it should be possible to specify how packets are processed and forwarded in a high-level domain-specific language (DSL) such as P4, and compiled to run on a software switch. We present PISCES, a software switch derived from Open vSwitch (OVS), a hard-wired hypervisor switch, whose behavior is customized using P4. PISCES is not hard-wired to specific protocols; this independence makes it easy to add new features. We also show how the compiler can analyze the high-level specification to optimize forwarding performance. Our evaluation shows that PISCES performs comparably to OVS and that PISCES programs are about 40 times shorter than equivalent changes to OVS source code.

Categories and Subject Descriptors: C.2.1 [Computer-Communication Networks] *Network Architecture and Design*; D.2.8 [Software Engineering] *Metrics—Complexity Measures; Performance Measures*

General Terms: Design; Languages; Performance

Keywords: Software-Defined Networks (SDN); Domain-Specific Languages (DSL); P4; Software Switch; OVS; Programmable Data Planes; PISCES; Compiler Optimizations

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

SIGCOMM'16, August 22–26, 2016, Florianópolis, Brazil.

Copyright 2016 ACM. ISBN 978-1-4503-4193-6/16/08...\$15.00

DOI: <http://dx.doi.org/10.1145/2934872.2934886>

1 Introduction

Software switches, such as Open vSwitch (OVS) [57], play a key role in modern data centers: with few exceptions, every packet that passes to or from a virtual machine (VM) passes through a software switch. In addition, servers greatly outnumber physical switches in this environment. Therefore, a data center full of servers running hypervisor software also contains far more software switches than hardware switches. Likewise, because each hypervisor hosts several VMs, such a data center has more virtual Ethernet ports than physical ones.

One of the main advantages of a software hypervisor switch is that it can be upgraded more easily than a hardware switch. As a result, hypervisor switches support new encapsulation headers, improved troubleshooting and debugging features, and middlebox-like functions such as load balancing, address virtualization, and encryption. In the future, as data center owners customize and optimize their infrastructure, they will continue to add features to hypervisor switches.

Each new feature requires customizing the hypervisor switch, yet making these customizations is more difficult than it may appear. First, most of the machinery that enables fast packet forwarding resides in the kernel. Writing kernel code requires domain expertise that most network operators lack, and thus introduces a significant barrier for developing and deploying new features. Recent technologies can accelerate packet forwarding in user space (*e.g.*, DPDK [34] and Netmap [64]), but these technologies still require significant software development expertise and intimate familiarity with a large, intricate, and complex codebase. Furthermore, customization requires not only incorporating changes into switch code, but also *maintaining* these customizations as the underlying software evolves over time, which can require significant resources.

Changing how a software switch forwards packets should not require intimate knowledge of how the switch is implemented. Rather, it should be possible to specify custom network protocols in a domain-specific language (DSL) such as P4 [10], which is then compiled to custom code for the hypervisor switch. Such a DSL would support customizing the forwarding behavior of the switch, without requiring changes to the underlying switch implementation. Decoupling custom protocol implementations from underlying switch code also makes it easier to maintain these customizations, since they

remain independent of the underlying switch implementation. With a standardized DSL, customizations may also be ported to other hardware or software switches, that support the same language.

A key insight, borrowed from a similar trend in hardware switches [11, 41], is that the underlying switch should be a substrate, well-tuned to process packets at high speed, but not tied to a specific protocol. In the extreme, the switch is said to be “protocol independent,” meaning that before it receives instructions about how to process packets (via a DSL), it does not know what a protocol is. Put another way, protocols are represented by programs written in the DSL, which protocol authors create.

We apply a similar philosophy to software switches. We assume the program written in the DSL specifies which packet headers to parse and the structure of the match-action tables (*i.e.*, which header fields to match and which actions to perform on matching headers). The underlying software substrate is a generic engine, optimized to parse, match, and act upon the packet headers in the form the program specifies.

Expressing these customizations in a DSL, however, entails compilation from the DSL to code that runs in the switch. Compared to a switch that is handwritten to implement fixed protocols, this protocol compilation process may reduce the efficiency of the underlying implementation and thus come at the cost of performance. The compilation process differs from hardware switches where, given limited resources, the objective is to optimize for metrics like area, latency, and power, while satisfying resource constraints [36]. Our goals in this paper are to (1) quantify the additional cost that expressing custom protocols in such a DSL produces; and (2) design and evaluate domain-specific compiler optimizations that reduce the performance overhead as much as possible. Ultimately, we demonstrate that, with the appropriate compiler optimizations, the performance of a *protocol-independent* software switch—a switch that supports custom protocol specification in a high-level DSL without direct modifications to the low-level source code—approaches parity with the native hypervisor software switch. Our results are promising, particularly given that OVS, our base code, was not designed to support protocol independence. Nevertheless, our results demonstrate that the “cost of programmability” in hypervisor switches is negligible. We expect our results will inspire the design of new protocol-independent software switches running at even higher speeds.

We make the following contributions:

- The design and implementation of PISCES, the first software switch that allows custom protocol specification in a high-level DSL, without requiring direct modifications to switch source code (Section 4).
- A public, open-source implementation of PISCES on GitHub [2]. The implementation is a protocol-independent software switch derived from OVS that is programmed from a high-level DSL, called P4.
- Domain-specific optimizations and a back-end optimizer to reduce the performance overhead of customizing OVS

using P4. We also introduce two new annotations in P4 to aid in the optimizations (Section 4.3).

- An evaluation of the code complexity of PISCES programs and its forwarding performance (Section 5). Our evaluation shows that PISCES programs are on average about 40 times shorter than equivalent changes to OVS source code and incur a forwarding performance (*i.e.*, throughput) overhead of only about 2%.

We begin by motivating the need for a customizable hypervisor software switch with a description of real use cases from operational networks (Section 2) and present background information on both P4 and OVS (Section 3).

2 The Need for a Protocol-Independent Switch

We say that PISCES is a *protocol-independent* software switch because it does not know what a protocol is or how to process packets on behalf of a protocol, until the programmer specifies it. For example, if we want PISCES to process IPv4 packets, then we need to describe how IPv4 packets are processed in a P4 program. In a P4 program (*e.g.*, `IPv4.p4`), we need to describe the format and fields of the IPv4 header, including the IP addresses, protocol ID, TTL, checksum, flags, and so forth. We also need to specify that we use a lookup table to store IPv4 prefixes, and that we search for the longest matching prefix. We also need to describe how a TTL is decremented, a checksum is updated, and so on. The P4 program captures the entire packet processing pipeline, which is compiled to source code for OVS that specifies the switch’s match, action, and parse capabilities.

A protocol-independent switch brings many benefits:

Adding new standard or private protocol headers. Vendors propose new protocol headers all the time, particularly for data centers. In recent years, VXLAN [47], NVGRE [73], Geneve [29] have all been standardized, and STT [16] and NSH [60] are also being discussed as potential standards. Private, proprietary protocols are also added, to provide a competitive advantage by, for example, creating better isolation between applications, or by introducing novel congestion marking. In many cases, before new protocols can be deployed, all hardware and software switches must be upgraded to recognize the headers and process them correctly. For hardware switches, the data center owner must provide requirements to their chip vendor and wait three to four years for the new feature to arrive, if the vendor agrees to add the feature at all. In the case of software switches, they must wait for the next major revision, testing, and deployment cycle. Even modifying an open-source software switch is not a panacea because once the data center owner directly modifies the open-source software switches to add their own custom protocols, these modifications still need to be maintained and synchronized with the mainline codebase, introducing significant code maintenance overhead as the original open-source switch continues to evolve. A data-center owner who could add new protocols

to a P4 program could, instead, compile and deploy a new protocol more quickly.

Removing a standard protocol header. Data-center networks typically run *fewer* protocols than legacy campus and enterprise networks, in part because most of the traffic is machine-to-machine and many legacy protocols are not needed (e.g., multicast, RSVP, L2-learning). For example, Amazon Web Services (AWS) reportedly only forwards packets using IPv4 headers [55]. It therefore benefits the data-center owner to remove unused protocols entirely, thus eliminating any concern of interactions with dormant implementations of legacy protocols. It is bad enough to have to support many protocols; much worse to have to understand interactions with and implications of protocols that operators do not intend to use. Therefore, data-center owners frequently want to eliminate unused protocols from their switches, NICs, and operating systems. Removing protocols from conventional switches is difficult; for hardware, it means waiting for new silicon, and for software switches it means wrestling with a large codebase to extract a specific protocol. In PISCES, removing an unused protocol is as simple as removing unused portions of a protocol specification and recompiling the switch source code. (Section 5.2.2 shows how this can even improve performance.)

Adding better visibility. As data centers get larger and are used by more applications, it becomes important to understand the network’s behavior and operating conditions. Failures can lead to huge loss in revenue, exacerbated by long debugging times as the network gets bigger and more complicated. There is growing interest in making it easier to see what the network is doing. Improving network visibility might entail supporting new statistics, generating new probe packets, or adding new protocols and actions to collect switch state (as is enabled by in-band network telemetry [42, 43]). Users will want to see how queues are evolving, latencies are varying, whether tunnels are correctly terminated, and whether links are still up. Often, during an emergency, users want to quickly add visibility features. Having them ready to deploy, or being able to modify forwarding and monitoring logic quickly may reduce the time to diagnose and fix a network outage.

Adding entirely new features. If users and network owners can modify the forwarding behavior, they may even add entirely new features. For example, over time we can expect switches to take on more complex routing, such as path-utilization aware routing [4, 40], new congestion control mechanisms [8, 19, 39], source-controlled routing [58], new load-balancing algorithms [26], new methods to mitigate DDoS [5, 25], and new virtual-to-physical gateway functions [17]. If a network owner can upgrade infrastructure to achieve greater utilization or more control, then they will know best how to do it. Given the means to upgrade a program written in a DSL like P4 for adding new features to a switch, we can expect network owners to improve their networks much more rapidly.



Figure 1: P4 abstract forwarding model.

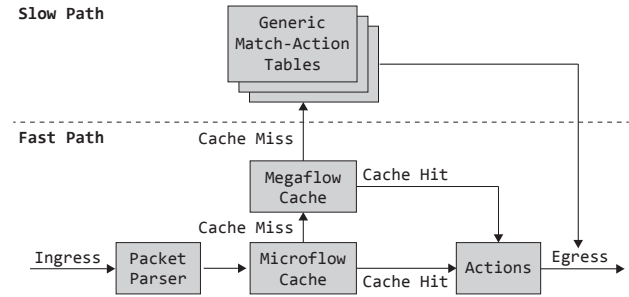


Figure 2: OVS forwarding model.

3 Background

PISCES is a *software switch* whose forwarding behavior is specified using a *domain-specific language*. PISCES is based on the Open vSwitch (OVS) [57] software switch and is configured using the P4 domain-specific language [10]. We describe both P4 and OVS below.

Domain-Specific Language: P4. P4 is a domain-specific language that expresses how the pipeline of a network forwarding element should process packets using the abstract forwarding model shown in Figure 1. In this model, each packet first passes through a programmable *parser*, which extracts headers. The P4 program specifies the structure of each possible header as well as a parse graph that expresses ordering and dependencies. Then, the packet passes through a series of match-action tables (MATs). The P4 program specifies the fields that each of these MATs may match and the control flow among them, as well as the spectrum of permissible actions for each table. At “runtime” (i.e., while the switch is forwarding packets), controller software may add, remove, and modify table entries with particular match-action rules that conform to the P4 program’s specification. Finally, a *deparser* writes the header fields back onto the packet before sending it out the appropriate port.

We choose P4 because its abstract model of a switch is similar to that of OpenFlow, the language built into OVS, which allows us to make straightforward apples-to-apples comparisons of OVS with and without a P4 front end. We considered other alternative bases, such as Click [44]—used in the Berkeley Extensible Software Switch (BESS) [30]—that allow for richer computation than match-action processing. However, for our purposes, P4 is sufficient to make the intended comparisons. There is merit to having a common way to express forwarding across all “plumbing” switches in a network, and have code that is portable from one to another. Therefore, using the same language makes sense for these experiments.

As BESS shows, there are other more extensible applications for software switches that are outside the scope of our work.

Software Switch: Open vSwitch. Open vSwitch (OVS) is widely used in data centers as a software switch running inside the hypervisor. In such an environment, OVS switches packets among virtual interfaces to VMs and physical interfaces. OVS implements common protocols such as Ethernet, GRE, and IPv4, as well as newer protocols found in data centers, such as the VXLAN Group Based Policy (GBP) extension [67], Geneve [29], NVGRE [73], and STT [16] for virtual network overlays.

The Open vSwitch virtual switch has two important pieces, called the *slow path* and the *fast path* (i.e., *datapath*), as shown in Figure 2. The slow path is a userspace program; it supplies most of the intelligence of OVS. The fast path acts as a caching layer that contains only the code needed to achieve maximum performance. Notably, the fast path must pass any packet that results in a cache miss to the slow path to get instructions for further processing. OVS includes a single, portable slow path and multiple fast-path implementations for different environments: one based on a Linux kernel module, another based on a Windows kernel module, and another based on Intel DPDK [34] userspace forwarding. The DPDK fast path yields the highest performance, so we use it for our work; with additional effort, our work could be extended to the other fast paths.

As an SDN switch, OVS relies on instructions from a controller to determine its behavior, specifically using the OpenFlow protocol [50]. OpenFlow specifies behavior in terms of a collection of *match-action tables*, each of which contains a number of entries called *flows*. In turn, a flow consists of a *match*, in terms of packet headers and metadata, *actions* that instruct the switch what to do when the match evaluates to true, and a numerical *priority*. When a packet arrives at a particular match-action table, the switch finds a matching flow and executes its actions; if more than one flow matches the packet, then the flow with the highest *priority* takes precedence.

A software switch that implements the behavior exactly as described above cannot achieve high performance, because OpenFlow packets often pass through several match-action tables, each of which requires general-purpose packet classification. Thus, OVS relies on caches to achieve good forwarding performance. The primary OVS cache is its *megaflow cache*, which is structured much like an OpenFlow [50] table. The idea behind the megaflow cache is that one could, in theory, combine all of the match-action tables that a packet visits while traversing the OpenFlow pipeline into a single table by computing their cross-product. This is infeasible, however, because the cross-product of k tables with n_1, \dots, n_k rules might have as many as $n_1 \times \dots \times n_k$ rules. The megaflow cache functions somewhat like a lazily computed cross-product: when a packet arrives that does not match any existing megaflow cache entry, the slow path computes a new entry, which corresponds to one row in the theoretical cross-product, and inserts

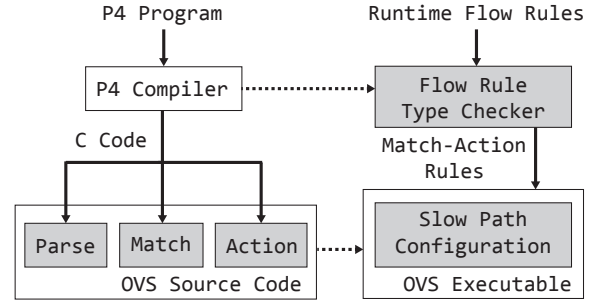


Figure 3: The P4-to-OVS Compiler in PISCES.

it into the cache. OVS uses a number of techniques to improve megaflow cache performance and hit rate [57].

When a packet hits in the megaflow cache, the switch can process it significantly faster than the round trip from the fast path to the slow path that a cache miss would require. As a general-purpose packet classification step, however, a megaflow cache lookup still has a significant cost. Thus, Open vSwitch fast-path implementations also include a *microflow cache*, a hash table that maps from a packet five-tuple to a megaflow cache entry. The result of the microflow cache lookup can only be a hint, because megaflows often match on more fields than just the five-tuple, so that a microflow cache entry can at best point to the most likely match. Thus, the fast path must verify that the megaflow cache entry indeed matches the packet. If it does match, the lookup cost is just that of the single hash table lookup. This lookup cost is generally much cheaper than general packet classification, so it is a significant optimization for traffic patterns with relatively long, steady streams of packets. If it does not match, then the packet continues through the usual megaflow cache lookup process, skipping the entry that it has already checked.

4 PISCES Prototype

Our PISCES prototype is a modified version of OVS with the *parse*, *match*, and *action* code replaced by C code generated by our P4 compiler. The workflow is as follows: First, the programmer creates a P4 program and uses the PISCES version of the P4 compiler (Section 4.1) to generate new *parse*, *match*, and *action* code for OVS. Second, OVS is compiled (using the regular C compiler) to create a protocol-dependent switch that processes packets as described in the P4 program. To modify a protocol, a user modifies the P4 program, which compiles to a new hypervisor switch binary.

We use OVS as the basis for PISCES because it is widely used and contains some basic scaffolding for a programmable switch, thus allowing us to focus only on the parts of the switch that need to be customized (i.e., *parse*, *match*, and *action*). The code is well-structured, lending itself to modification, and test environments already exist. It also allows for apples-to-apples comparisons: We can compare the number of lines of code in unmodified OVS to the P4 program for PISCES (Section 5.1), and we can also compare their performance (Section 5.2).

4.1 The P4-to-OVS Compiler in PISCES

P4 compilers have two parts: a front end that turns the P4 code into a target-independent intermediate representation (IR), and a back end that maps the IR to the target. In our case, the back end optimizes CPU time, latency, or other objectives by manipulating the IR, and then generates C code that replaces the parsing, match, and action code in OVS, as shown in Figure 3. The P4-to-OVS compiler outputs C source code that implements everything needed to compile the corresponding switch executable. The compilation process also generates an independent type checking program that the executable uses to ensure that any runtime configuration directives from the controller (*e.g.*, insertion of flow rules) conforms to the protocol specified in the P4 program.

Parse. The C code that replaces the original OVS *parser* is created by replacing `struct flow`, the C structure that OVS uses to track protocol header fields, to include a member for each field specified by the P4 program, and generating code to extract header fields from a packet into `struct flow`.

Match. OVS uses a general-purpose classifier data structure, based on tuple-space search [69], to implement matching. To perform custom matches, we do not need to modify this data structure or the code that manages it. Rather, the control plane can simply populate the classifier with new packet header fields at runtime, thereby automatically making those fields available for packet matching.

Action. The back end of our compiler supports custom actions by automatically generating code that we statically compile into the OVS binary. Custom actions can execute either in the OVS slow path or the fast path; the compiler determines where a particular action will run to ensure that the switch performs the actions efficiently. Certain actions (*e.g.*, `set_field`) can execute in either component. The programmer can offer hints to the compiler as to whether slow path or fast path implementation of an action is most appropriate.

Control flow. In a switch, a packet’s *control flow* is the sequence of match-action tables that the packet traverses. Whereas with P4, control flow must be specified at the program’s compile time, in OVS control flow is specified at runtime, via flow entries, which makes it more flexible. Therefore, our compiler back end can implement P4 control semantics without OVS changes.

Optimizing the IR. The compiler back end contains an optimizer to examine and modify the IR, so as to generate high-performance C code. For example, a P4 program may include a complete IP checksum, but the optimizer can turn this operation into an incremental IP checksum to make it faster. The compiler also performs data-flow analysis on the IR [3], allowing it to coalesce and specialize the C code. The optimizer also decides when and where in the packet processing pipeline to edit packet headers. Some hardware switches postpone editing until the end of the pipeline, whereas software switches typically edit headers at each stage in the pipeline. If necessary,

the optimizer converts the IR for in-line editing. We describe the optimizer in more detail in Section 4.3.

As is the case with other P4 compilers [10, 36], the P4-to-OVS compiler also generates an API for the match-action tables, and extends the OVS command-line tools to work with the new fields.

4.2 Modifications to OVS

We need to make three modifications to OVS to enable it to implement the forwarding behavior described in any P4 program.

Arbitrary encapsulation and decapsulation. OVS does not support arbitrary encapsulation and decapsulation, which a P4 program might require. Each OVS fast path provides custom support for various fixed forms of encapsulation. The Linux kernel fast path and DPDK fast path, for example, each separately implement GRE [22], VXLAN [47], STT [16], and other encapsulations. The metadata required to encapsulate and decapsulate a packet for a tunnel is statically configured. The switch uses a packet’s ingress port to map it to the appropriate tunnel; on egress, the packet is encapsulated in the corresponding IP header based on this static tunnel configuration. We therefore added two new primitives to OVS, `add_header()` and `remove_header()`, to perform encapsulation and decapsulation, respectively, and perform these operations in the fast path.

Conditionals based on comparison of header fields. OpenFlow directly supports only bitwise equality tests against header fields. Relational tests such as `<` and `>` to compare a k -bit field against a constant can be expressed as at most k rules that use bitwise equality matches. A relational test between two k -bit fields, such as $x < y$, requires $k(k+1)/2$ such rules. To simultaneously test for two such conditions that individually take n_1 and n_2 rules, one needs $n_1 \times n_2$ rules. P4 directly supports such tests, but implementing them in OpenFlow this way is too expensive, so we added direct support for them in OVS as *conditional actions*, a kind of “if” statement for OpenFlow actions. For example, our extension allows the P4 compiler to emit an action of the form “If $x < y$, go to table 2, otherwise go to table 3.”

General checksum verify/update. An IP router should verify the checksum at ingress, and recompute it at egress, and most hardware switches do it this way. A software router often skips checksum verification on ingress to reduce CPU cycles. Instead, it just incrementally updates the checksum if it changes any fields (*e.g.*, the TTL).¹ Currently, OVS only supports incremental checksums, but we want to support other uses of checksums in the way the programmer intended. We therefore added incremental checksum optimization, described in Section 4.3. Whether this optimization is valid depends on whether the P4 switch is acting as a forwarding element or an end host for a given packet—if it is an end host, then it

¹If the checksum was incorrect before the update, it is still incorrect afterward, and we rely on the ultimate end host to discard the packet.

Optimization	CPU Cycles	Slow-Path Trips
Inline- vs. post-pipeline editing	✓	
Incremental checksum	✓	
Parser specialization	✓	
Action specialization	✓	
Action coalescing	✓	
Cached field modifications	✓	✓
Stage assignment	✓	✓

Table 1: Back-end optimizations and how they improve performance.

must verify the checksum—so it requires annotation by the P4 programmer.

4.3 The Compiler’s Back-end Optimizer

Two aspects of a software switch ultimately affect forwarding performance: (1) the per-packet cost for fast-path processing (adding 100 cycles to this cost reduces the switch’s throughput by about 500 Mbps), and (2) the number of packets sent to the slow path, which takes 50+ times as many cycles as the fast path to process a packet. Table 1 lists the optimizations that we have implemented, as well as whether the optimization reduces trips to the slow path, fast path CPU cycles, or both. The rest of the section details these optimizations.

Inline editing vs. post-pipeline editing. The OVS fast path performs *inline editing*, applying packet modifications immediately (the slow path does some simple optimization to avoid redundant or unnecessary modifications). If many header fields are modified, removed or inserted, it can become costly to move and resize packet data on the fly. Instead, it can be more efficient to delay editing until the headers have been processed (as hardware switches typically do). The optimizer analyzes the IR to determine how many times a packet may need to be modified in the pipeline. If the value is below a certain threshold, then the optimizer performs inline editing; otherwise, it performs post-pipeline editing. We allow the programmer to override this heuristic using a pragma directive.

Incremental checksum. By expressing a checksum operation in terms of a high-level program description such as P4, a programmer can provide a compiler with the necessary contextual information to implement the checksum more efficiently. For example, the programmer can inform the compiler via annotations that the checksum for each packet can be computed incrementally [51]; the optimizer can then perform data-flow analysis to determine which packet header fields change, thus making re-computation of the checksum more efficient.

Parser specialization. Protocol-independent software switches can optimize the implementation of the packet parser, since a customized packet processing pipeline (as specified in a high-level language such as P4) provides specific information about which fields in the packet are modified or used as the basis for forwarding decisions. For example, a layer-2 switch that does not make forwarding decisions based on information at other layers can avoid parsing packet header fields at those layers. Specifying the forwarding behavior in a

high-level language provides the compiler with information that it can use to optimize the parser.

Action specialization. The inline editing actions in the OVS fast path group together related fields that are often set at the same time. For example, OVS implements a single fast path action that sets the IPv4 source, destination, type of service, and TTL value. This is efficient when more than one of these fields is to be updated at the same time, with little marginal cost if only one is updated. IPv4 has many other fields, but the fast path cannot set any of them.

The design of this aspect of OVS required domain expertise: its designers knew which fields were important for the fast path to be able to change. A P4 compiler does not have this kind of expert knowledge of which fields to group together, yielding a possible cost for grouping too few or too many fields into a single action. Fortunately, the high-level P4 description of the match-action control flow allows the optimizer to identify and eliminate redundant checks in the fast-path set actions, using optimizations like dead-code elimination [3]. This way, the optimizer only checks those fields in the set actions that will actually be set in the match-action control flow.

Action coalescing. By analyzing the control flow and match-action processing in the P4 program, the compiler can discover which fields are actually modified and can generate an efficient, single action to directly update those fields. Thus, if a rule modifies two fields, the optimizer only installs one action in OVS.

Cached field modifications. Network protocol data planes rarely require arithmetic operations on header fields. TTL decrement operations are the most obvious counterexample; checksums, already addressed above, are another. Thus, OVS fast paths do not include general-purpose arithmetic operations. In fact, they do not include a special-purpose TTL decrement operation either. Instead, to implement the special-purpose OpenFlow action to decrement a TTL, the slow path relies on the fact that most packets from a given source have the same TTL. Therefore, it emits a cache entry that matches on the TTL value observed in the packet that it is forwarding and overwrites this value with one less than that observed value, an approach we call “match-and-set.” For TTL decrement, this solution is acceptable because the OVS designers know that caching this way yields a high hit rate in practice.²

Match-and-set is not always appropriate. As a straw man, consider update of the IPv4 or IPv6 checksum given a change in some other IP field. With a match-and-set approach, the cache entry would have to match on every field that contributes to the checksum, that is, every IP field, which would reduce the cache entry’s hit rate nearly to zero. The same can be true for simpler arithmetic operations that P4 supports, such as incrementing or decrementing a field value, and in the end PISCES has no way to know whether match-and-set is appropriate in a given case.

²In addition, real-world uses of TTL decrement are always paired with a “TTL exceeded” check that would itself cause the cache entry to match on TTL, which would negate the value of a special-case TTL decrement action.

The solution that PISCES takes is to avoid match-and-set when it can, by automatically generating fast path operations to implement the particular arithmetic operations that a P4 program requires. For example, if the program increments a particular field, PISCES generates a fast path operation to increment that field. This is effective when the P4 program executes the arithmetic operation “blindly,” without otherwise matching on the modified field’s value. If the program does match on it, then, following the usual rules for caching, the cache entry must match on the field, so that a match-and-set approach is necessary.

Stage assignment. OVS implements *staged lookup* [57] to reduce the number of trips to the slow path. Staged lookup divides fields into a ordered list of groups, called *stages*. The stages are cumulative, so that each stage after the first contains all of the fields from the previous stages plus additional fields. The final stage contains every field. OVS implements each stage as a separate hash table in its tuple space search classifier. A classifier lookup searches each of these stages in order. If any search yields no match, the overall search terminates and only the fields included in the last stage must be matched in the cache entry.

OVS uses four such stages: the first stage is metadata fields (such as the packet’s ingress port), the second is metadata and layer-2 fields, the third adds layer-3 fields, and the fourth includes all fields (*i.e.*, metadata, layer 2, 3, and 4). This order is based on the principle that stages are most effective when their order corresponds to increasing order of entropy in the observed values of fields for networks [66]. In the common case, for example, a cache entry that matches on metadata only is likely to have a higher hit rate than a cache entry that matches only on layer-4 fields, so metadata first appears in an earlier stage (the first stage) than do layer-4 fields (the final stage).

Staged lookup generalizes to arbitrary P4 programs. This ordering cannot be inferred from the P4 program, so PISCES needs assistance to choose appropriate stages. We augmented the P4 language to enable a user to annotate each header with a stage number. The number of stages is the same as the number of headers.

5 Evaluation

We compare the complexity and performance of a PISCES virtual software switch with equivalent OVS native packet processing. We compare the resulting programs along two dimensions: (1) complexity, including development and deployment complexity as well as maintainability; (2) performance, by comparing packet-forwarding performance of PISCES to the same native OVS functionality.

5.1 Complexity

Complexity indicates the ease with which a program may be modified to fix defects, meet new requirements, simplify future maintenance, or cope with changes in the software environment. We evaluate two categories of complexity: (1) *de-*

	LoC	Methods	Method Size
OVS	14,535	106	137.13
PISCES	341	40	8.53

Table 2: Native OVS compared to equivalent baseline functionality implemented in PISCES.

		Files Changed	Lines Changed
Connection Label:	OVS [70, 71]	36	633
	PISCES	1	5
Tunnel OAM Flag:	OVS [27, 28]	21	199
	PISCES	1	6
TCP Flags:	OVS [61]	20	370
	PISCES	1	4

Table 3: The number of files and lines we needed to change to implement various functionality in P4, compiled with PISCES, compared to adding the same functionality to native OVS.

velopment complexity of developing baseline features for a software switch; and (2) *change complexity* of maintaining an existing software switch.

5.1.1 Development complexity

We evaluate development complexity with three different metrics: lines of code, method count, and average method size. We count lines of code simply by counting line break characters and the number of methods by counting the number of subroutines in each program, as measured using `ctags` [33]. Finally, we divide lines of code by number of methods to arrive at the average method size. A high average might indicate that (some) methods are too verbose or complex.

Writing a compiler is a one-time cost. Whereas developers update their P4 programs frequently, the compiler is changed much less often—usually when the P4 language specification changes. For PISCES, we write about 1,000 lines of code for compiling P4 to C code, and an extra 1,700 lines of code to extend the native OVS to incorporate the generated C code.

`ovs.p4` [1] contains the representation of the headers, parsers, and actions that are currently supported in OVS. Much of the code in OVS is out of the scope of P4, so our measurements include only the files that are responsible for protocol definitions and header parsing. Table 2 summarizes each of these metrics for the native OVS header fields and parser implementation, and the equivalent logic in P4.³ PISCES reduces the lines of code by about a factor of 40 and the average method size by about a factor of 20.

5.1.2 Change complexity

To evaluate the complexity of *maintaining* a protocol-independent software switch in PISCES, we compare the effort

³We reuse the same code for the match-action tables in both implementations because this logic generalizes for both OVS and a protocol-independent switch such as PISCES.

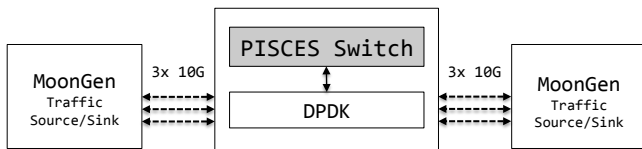


Figure 4: Topology of our evaluation platform.

required to add support for a new header field in a protocol that is otherwise already supported, in OVS and in P4. Table 3 shows our analysis of changes to add support for three fields: (1) *connection label*, a 128-bit custom metadata to the connection tracking interface; (2) *tunnel OAM flag*, which many networking tools use to distinguish test packets from real traffic; and (3) *TCP flags*, a modification that adds support for parsing all of the TCP flags. Table 3 shows the changes to OVS based on the public Open vSwitch commits. These numbers are conservative because they include only the changes to one of the three OVS fast-path implementations.

The results demonstrate that modifying just a few lines of code in a single P4 file is sufficient to support a new field, whereas in OVS, the corresponding change often requires hundreds of lines of changes over tens of files. Among other changes, one must add the field to `struct flow`, describe properties of the field in a global table, implement a parser for the field in the slow path, and separately implement a parser in one or more of the fast paths.

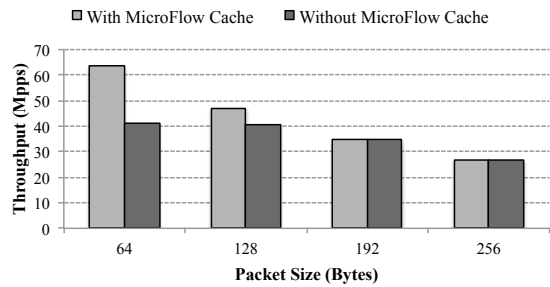
5.2 Forwarding Performance

In this section, we compare OVS and PISCES packet-forwarding performance.

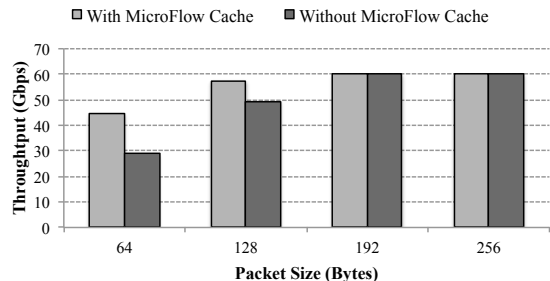
5.2.1 Experiment setup and evaluation metrics

Figure 4 shows the topology of the setup for evaluating the forwarding performance of PISCES. We use three PowerEdge R730xd servers with two 8-core, 16-thread Intel Xeon E5-2640 v3 2.6GHz CPUs running the Proxmox Virtual Environment [59], an open-source server virtualization platform that uses virtual switches to connect VMs, with Proxmox Kernel version 4.2.6-1-pve. Each of our machines is equipped with one dual-port and one quad-port Intel X710 10 Gbps NIC. We configure two such machines with MoonGen [20] to send minimum-size 64-byte frames at 14.88 million packets per second (Mpps) full line rate on three of the 10 Gbps interfaces [64], leaving the other interfaces unused. We connect these six interfaces to a third machine, the device under test, sending a total of 60 Gbps of traffic for PISCES to forward.

We consider throughput and packets-per-second to compare the forwarding performance of PISCES and OVS, using the MoonGen packet generator to generate test traffic for our experiments. We configure PISCES and OVS with six Poll Mode Driver (PMD) threads—one for each 10 Gbps interface—in a Run-to-Completion (RTC) model [35]. Each thread runs on a separate CPU core attached to one of the Non-Uniform Memory Access (NUMA) [45] nodes on the machine. To further understand performance bottlenecks, we use the machine’s time-stamp counter (TSC) to measure the number of



(a) Forwarding performance in millions of packets per second with a standard deviation of less than 0.035 Mpps for all data points.



(b) Forwarding performance in gigabits per second with a standard deviation of less than 0.026 Gbps for all data points.

Figure 5: Forwarding performance for OVS with and without the microflow cache enabled, for input traffic of 60 Gbps across all six ports and one flow rule per port.

CPU cycles used by various packet processing operations (*i.e.*, parser, megaflow cache lookup, and actions). When reporting CPU cycles, we report the average CPU cycles per packet over all packets forwarded in an experiment run; each run lasts for 30 seconds and has an ingress rate of 89.28 Mpps.

Calibrating OVS to enable performance comparison. To more accurately measure the cost of parsing for both OVS and PISCES in subsequent experiments, we begin by establishing a baseline for OVS performance with minimal parsing functionality. To minimize the cost of parsing, we disable the parser, which ordinarily parses a comprehensive fixed set of headers, so that it reports only the input port. After this change, we send test traffic through the switch with a trivial flow table that matches every packet that ingresses on port 1 and sends it to port 2.

We measured the performance of this modified OVS. Figures 5a and 5b show the maximum throughput that our setup achieves with OVS, with and without the microflow cache, for 60-Gbps traffic. For 64-byte packets, disabling the microflow cache reduces performance by about 35%, because a lookup in the OVS megaflow cache consumes five times as many cycles as the microflow cache (Table 4). For small packets, the OVS switch is CPU-bound on lookups; thus, in this operating regime, the benefit of the microflow cache is clear.

With this calibration in mind, for the remainder of this section, we use the forwarding performance for OVS with the

Switch Components	With MicroFlow	Without MicroFlow
Parser	19.0	18.9
MicroFlow Cache	18.9	—
MegaFlow Cache	—	92.2
Slow Path	—	—
Fast-Path Actions	39.9	38.8
End-to-End	100.6	166.0

Table 4: Average number of cycles per packet consumed by each element in the virtual switch when processing a 64-byte packet.

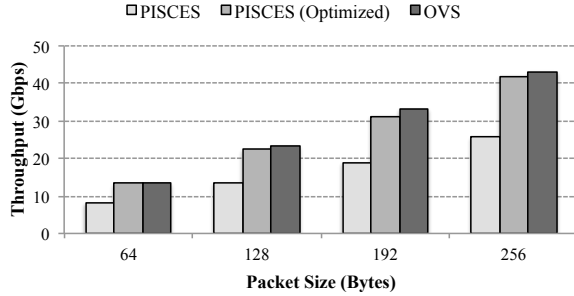


Figure 6: Throughput comparison of L2L3-ACL benchmark application between OVS and PISCES in gigabits per second, with a standard deviation of less than 0.023 Gbps for all data points.

microflow cache disabled as the basis for our performance comparison to PISCES. We disable the microflow cache because it relies on matching a hash of a packet’s five-tuple, which most NICs can compute directly in hardware. Although OVS’s microflow cache significantly improves its forwarding performance, this feature relies on protocol-dependent features (specifically, that the packet has a five-tuple in the first place). Because our goal is to evaluate forwarding rates for protocol-independent switches, we disabled OVS’s microflow cache so that we could compare PISCES, a protocol-independent switch, with a version of OVS that has no protocol-dependent optimizations. Comparing PISCES performance to that of OVS with microflow caching disabled thus offers a more apples-to-apples performance comparison, although it makes it difficult to interpret performance versus “real-life Open vSwitch.” We expect that implementing a microflow cache in PISCES, by adding P4 annotations for the fields to be hashed and then hashing them in software, would recover most of the performance.

5.2.2 End-to-end performance

We next measure the forwarding performance of a real-world network application for both OVS and PISCES. This evaluation provides a clear illustration of the end-to-end performance costs of programmability. We select a realistic and relatively complex application where both switch implementations provide all packet processing features to provide a fair performance comparison of PISCES in realistic network settings.

Figure 7 shows this application, which we call “L2L3-ACL.” It performs the following operations:

- Parse Ethernet, VLAN, IP, TCP and UDP protocols.
- Perform VLAN encapsulation and decapsulation.
- Perform control-flow and match-action operations according to Figure 7 to implement an access control list (ACL).
- Set Ethernet source, destination, type and VLAN fields.
- Decrement IP’s TTL value.
- Update IP checksum.

Table 5 shows the forwarding performance results for this application. The most important rows are the last two, which show a “bottom line” comparison between OVS and PISCES, after we apply all compiler optimizations. These results show that both the average number of CPU cycles per packet and the average throughput for PISCES with all compiler optimizations is comparable to OVS with microflow caching disabled: both require just over an average of 400 CPU cycles per packet, and both achieve throughput of just over 13 Gbps—a *performance overhead of less than 2%*. Figure 6 demonstrates that this result also holds for larger packet sizes. In all cases, PISCES with compiler optimizations enabled in its compiler achieves performance comparable to OVS.

Next, we discuss in more detail the performance benefits that each compiler optimization achieves for this end-to-end application.

Individual compiler optimizations. P4 supports post-pipeline editing, so we start by compiling L2L3-ACL with post-pipeline editing. PISCES requires an average of 737 cycles to process a 64-byte packet. Packet parsing and fast-path actions are primarily responsible for these additional CPU cycles. As our microbenchmarks demonstrate (Section 5.2.3), if the number of adjustments to packets are less than eight, using inline-editing mode provides better forwarding performance. Based on that insight, the PISCES version of the P4 compiler uses inline editing, which reduces the number of cycles consumed by the parser by about 56%. However, fast-path actions’ cycles slightly increased (still 255 cycles more than OVS).

Next, we introduce incremental checksum updates to reduce the number of cycles consumed by the fast-path actions. The only IP field that is modified is TTL, but the full checksum verify and update design supported by P4 abstract model runs the checksum over entire headers once at the ingress and once at egress. For our P4 program, we specify that we want to use incremental checksum. Using this knowledge, instead of recalculating checksum on all header fields, using data-flow analysis on the P4 program (MAT and control-flow), the P4 compiler determines that the pipeline modifies only the TTL and adjusts the checksum using only that field, which reduces the number of cycles consumed by the fast-path actions by 59.7%, a significant improvement. However, PISCES still consumes 23.24 more cycles than OVS.

To further improve the performance we apply action specialization and coalescing, and parser specialization (Section 4.3).

Switch	Optimization	Parser	MegaFlow Cache	Fast-Path Actions	End-to-End (Avg.)	Throughput (Mbps)
PISCES	Baseline	76.5	209.5	379.5	737.4	7590.7
	Inline Editing	-42.6	—	+7.5	-45.4	+281.0
	Inc. Checksum	—	—	-231.3	-234.5	+4685.3
	Action Specialization	—	—	-10.3	-9.2	+191.2
	Parser Specialization	-4.6	—	—	-7.6	+282.3
	Action Coalescing	—	—	-14.6	-14.8	+293.0
	All optimizations	29.7	209.0	147.6	425.8	13323.7
OVS	—	43.6	197.5	132.5	408.7	13497.5

Table 5: Improvement in average number of cycles per packet, consumed by each element in the virtual switch when processing 64-byte packet, for L2L3-ACL benchmark application. (Most listed optimizations for the PISCES version of the P4 compiler do not have any counterpart in OVS, but OVS does implement incremental checksums.)

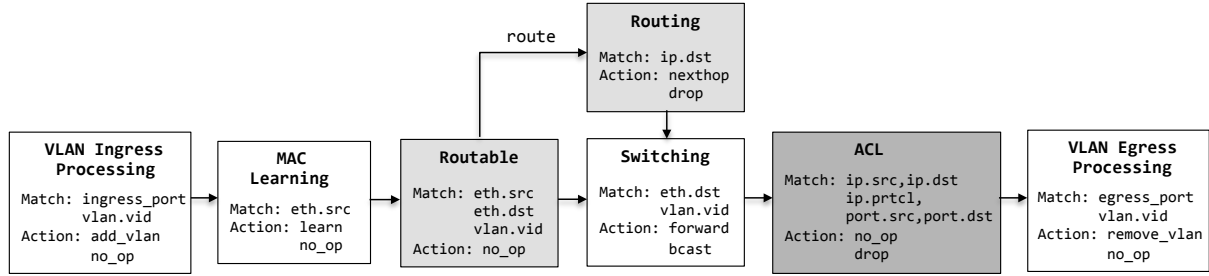


Figure 7: Control flow of L2L3-ACL benchmark application. Each of these tables contains a list of fields to match on and a set of actions to choose from when installing a flow rule. For example, in VLAN Ingress Processing, one can match on ingress port and VLAN id, and can perform add_vlan or no_op actions.

This brings the number of cycles consumed per packet by PISCES to 425.82.

Parser specialization. A protocol-independent switch only needs to parse the packet-header fields for the protocols defined by the programmer. The compiler in PISCES can optimize the parser further to only parse the header fields that the switch needs to process the packet. To evaluate the potential benefits of this specialization, we repeat our end-to-end performance evaluation using two subsets of the L2L3-ACL program: the “L2L3” program, which does not perform the ACL functions, and the “L2” program, which manipulates the Ethernet and VLAN headers and performs VLAN encapsulation, but which does not parse any IP headers or decrement the TTL (and thus does not update the IP checksum). In terms of the control flow from the original “L2L3-ACL” benchmark program from Figure 7, the “L2L3” program removes the dark grey ACL tables, and the “L2” program additionally removes the light grey Routable and Routing tables.

Table 6 compares the forwarding performance of OVS and PISCES for these two programs. For L2L3, PISCES consumes four more cycles per packet than OVS. However, PISCES has faster parsing: compared to L2L3-ACL, parsing in L2L3 is about seven cycles per packet cheaper. OVS uses a fixed parser, so its cost remains constant. Parser specialization removes redundant parsing of fields from the parser that are not used in the control-flow (*i.e.*, TCP and UDP headers). Because OVS does not know the control-flow and MAT structure a

priori, its parser cannot achieve the same specialization. In the case of the L2 application, the parser could specialize further, since it needs only to parse Ethernet headers. In this case, PISCES can actually process packets *more quickly* than the protocol-dependent switch.

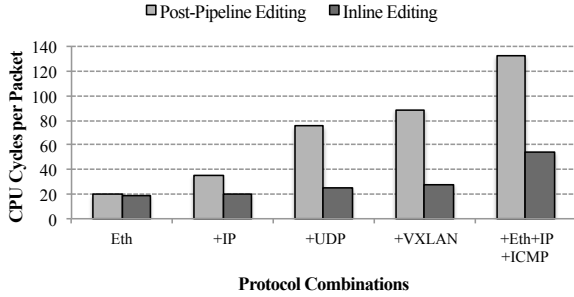
5.2.3 Microbenchmarks

We now evaluate the performance of individual components of PISCES. We focus on the parser and actions, which are applied on every incoming packet and have the largest effect on performance. We now benchmark how increasing complexity in both parser and actions affect the overall performance of PISCES.

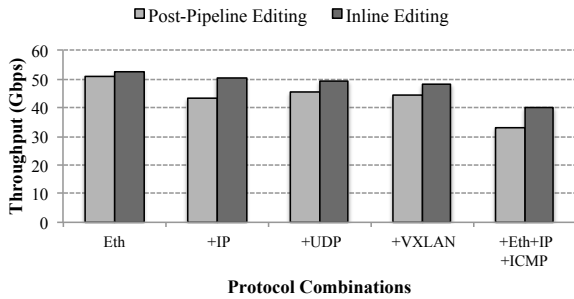
Parser performance. Figure 8a shows how per-packet cycle counts increase as the P4 program parses additional protocols, for both post- and inline-editing modes. To parse only the Ethernet header, the parser consumes about 20 cycles, in either mode. As we introduce new protocols, the cycle count increases, more rapidly for post-pipeline editing, for which the switch creates an extra copy of the protocol headers for fast-path actions. For the largest protocol combination in Figure 8a, the parser requires about 133 cycles (almost six times as many cycles as simply processing an Ethernet frame) for post-pipeline editing and 54 cycles for inline-editing. Figure 8b shows how the throughput decreases with the addition of each new protocol in the parser. For input traffic at 60 Gbps, switching throughput decreases about 35%, from 51.1 Gbps

Switch	Programs	Optimizations	Parser	MegaFlow Cache	Fast-Path Actions	End-to-End (Avg.)	Throughput (Mbps)
PISCES	L2L3	Optimized	22.9	188.4	130.5	392.3	14159.1
OVS	L2L3	—	43.6	176.0	131.8	388.3	14152.2
PISCES	L2	Optimized	19.7	148.2	90.9	305.7	18118.5
OVS	L2	—	43.6	155.2	78.7	312.1	17131.3

Table 6: Improvement in average number of cycles per packet, consumed by each element in the virtual switch when processing 64-byte packet, for L2L3 and L2 benchmark applications.



(a) CPU cycles.



(b) End-to-end throughput with a standard deviation of less than 0.063 Gbps for all data points.

Figure 8: Effect on parser CPU cycles and end-to-end throughput as more protocols are added to the parser.

to 33.2 Gbps, for post-pipeline editing and about 24%, from 52.4 Gbps to 40.0 Gbps, for inline editing.

Fast-path action performance. Performance-wise, the dominant action in a virtual switch is the set-field (or modify-field) action or, in other words, a write action. Figure 9 shows the per-packet cost, in cycles, as we increase the number of set-field actions in the fast path for both post- and inline-editing modes. In post-editing mode, we apply our changes to a copy of the header fields (extracted from the packet) and at the end of the pipeline execute a “deparsed” action that writes the changes back to the packet. The “deparsed” bar shows how deparsing consumes about 99 cycles even if no fields are modified, whereas inline editing has no cost in this case. As the number of writes increases, the performance difference between the two modes narrows. For 16 writes, this difference is 20 cycles less than for a single write. Still, in both cases, the number of cycles increases. For post-editing case, 16 writes consumes 354 cycles, about 3.6 times that of a single write;

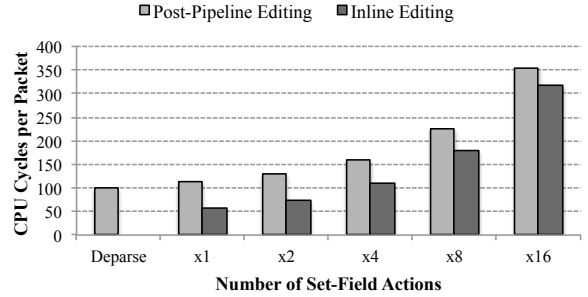


Figure 9: Fast Path Set-Field Action Performance.

for inline editing, 16 writes consumes 319 cycles, or about 5.6 times more cycles than a single write.

We also measured cycles-per-packet for adding or removing headers. Figures 10 and 11 show cycles-per-packet for an increasing number of add-header and remove-header actions, respectively, in the post-pipeline and inline-editing modes.

For the add-header action, for inline-editing mode, the number of cycles doubles for every new action. This is because these actions are applied directly on the packet, adjusting the packet size each time. In contrast, post-pipeline editing adjusts the packet size only once, in the “deparsed” action, so that the number of cycles consumed remains almost constant. For a single add-header action, post-editing cost is higher, but for four or more actions the inline-editing mode is more costly. For 16 add-header actions, inline editing consumes 577 more cycles per packet than post-pipeline editing.

We observe a similar trend for remove-header action. There is one additional wrinkle: as the number of remove-header actions increases, the cost of post-pipeline editing actually decreases slightly, because fewer bytes need to be adjusted in the packet as the packet shrinks. As we increase the number of remove-header actions from 1 to 16, the per-packet cycle count decreases by about 21%. This led us to the following rule of thumb: for fewer than 8 packet-size adjustments (*i.e.*, add- and remove-header actions), the compiler uses inline-editing; otherwise, it applies post-pipeline editing, as the added number of cycles required by the parser to generate a copy of the parsed packet headers is offset by the number of cycles required by the add/remove header actions in the inline-editing mode.

Slow-path forwarding performance. When OVS must send all packets to the slow path, it takes on average about 3,500 cycles to process a single packet (about 50 times the cycles incurred for a microflow cache hit). In this case, the maximum

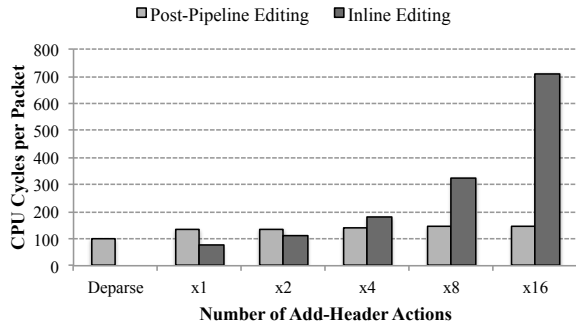


Figure 10: Fast Path Add-Header Performance.

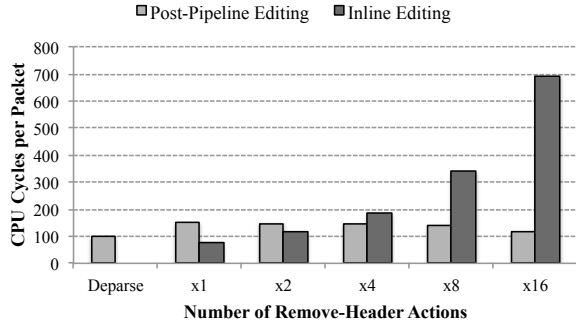


Figure 11: Fast Path Remove-Header Performance.

packet forwarding rate is about 0.66 Mpps regardless of packet size. This per-packet cycle count for slow-path processing was for the simplest possible program that sends every packet to the same output port. Most real packet processing programs would require significantly more cycles. For example, for the L2L3-ACL program, slow-path processing required anywhere from 30,000 to 60,000 CPU cycles per packet. These performance numbers indicate the importance of the megaflow cache optimizations that we described in Section 4.3 to reduce the number of trips to the slow path. Clearly, the number of trips to the slow path depends on the actual traffic mix (because this affects cache hit rates in the megaflow cache), so it is difficult to state general results about the benefits of these optimizations, but computing the slowdown as a result of cache misses is straightforward.

Control flow. Control flow in OVS, and thus in PISCES, is implemented in the slow path. It has a small one-time cost, which is impossible to separate from slow path performance in general, at the setup of every new flow.

6 Related Work

PISCES protocols and packet-processing functions can be specified using a high-level domain-specific language for packet processing. Although PISCES uses P4 as its high-level language and OVS as its software switch, previous work has developed both domain-specific languages for packet processing and virtual software switches, where our approaches

for achieving protocol independence and efficient compilation from a DSL to a software switch may also apply.

Domain-specific languages for packet processing. The P4 language provided the main framework for protocol independence [10]; PISCES realizes protocol independence in a real software switch. P4 itself borrows concepts from prior work [7, 23, 48]; as such, it may be possible to apply similar concepts that we have implemented in PISCES to other high-level languages. Although PISCES compiles P4 to OVS source code, the concepts and optimizations that we have developed could apply to other high-level languages and target switches; an intermediate representation such as NetASM [65] could ultimately provide a mechanism for a compiler to apply optimizations for a variety of languages and targets. Languages such as Pyretic [62] and Frenetic [24] are domain-specific languages that specify how packets should be processed by a fixed-function OpenFlow switch. They would require significant adaptation to take advantage of the abilities of a programmable switch. Also, compiling packet programs to reconfigurable hardware switches [36] and FPGAs [12, 63] differs from compiling to software switches. For hardware switches, the focus is on constrained optimization problems where, given a relatively small chip or memory footprint, the goal is to use that space optimally while satisfying dependencies. Such an approach is not likely to be effective for software switches, which do not have the same kinds of constraints.

Virtual software switches. Existing methods and frameworks for building software switches like Linux Kernel [46], DPDK [34], Netmap [64], Click [44], and BPF [14, 15, 49] require intimate knowledge about the underlying implementation and, thus, make it difficult for a network programmer to rapidly adapt and add new features to these virtual switches. PISCES, on the other hand, allows programmer to specify packet processing behavior independent of the underlying implementation details. Open vSwitch (OVS) [57] provides interfaces for populating its match-action tables but does not provide mechanisms to customize protocols and actions.

Other programmable switches. Software routers such as RouteBricks [18], PacketShader [31], and GSwitch [72] rely on general-purpose processors or GPUs to process packets; these designs generally focus on optimizing server, network interface, and processor scheduling to improve the performance of the software switch. These switches do not enable programmability through a high-level domain-specific language such as P4, and they also do not function as hypervisor switches. CuckooSwitch [74] can be used as a hypervisor switch. However, it focuses on providing fast forwarding table lookups by using highly-concurrent hash tables based on Cuckoo hashing [54], and, also does not provide a high-level domain-specific language to configure the switch. SwitchBlade [6] enables some amount of protocol customization and forwards packets at hardware speeds, but also acts as a standalone switch and requires an FPGA as a target.

Measuring performance. Previous work has both measured [9, 21] and improved [14, 15, 34, 49, 56, 57, 64] the perfor-

mance of software virtual switches. Work on measurement has converged on a set of performance metrics to compare various switch architectures and implementations; our evaluation uses these metrics to compare the performance of PISCES to that of other virtual switches.

Measuring complexity. A number of metrics for measuring the complexity and maintainability of a program written in a domain-specific language are developed in software engineering [13, 32, 37, 38, 52]. One of the goal of PISCES is to make it easier for the programmer to develop and maintain code. For our evaluation, we have taken these metrics from software engineering to evaluate the complexity of writing a program in P4 vs. directly modifying the OVS source code in C.

7 Conclusion

The increasing use of software hypervisor switches in data centers has introduced the need to rapidly modify the packet forwarding behavior of these software switches. Today, modifying these switches requires both intimate knowledge of the switch codebase *and* extensive expertise in network protocol design, making the bar for customizing these software switches prohibitively high. As an alternative to this mode of operation, we developed PISCES, a programmable, protocol-independent software switch that allows a protocol designer to specify a software switch’s custom packet processing behavior in a high-level domain-specific language (in our case, P4); a compiler then produces source code for the underlying target software switch (in our case, OVS). PISCES programs are about 40 times more concise than the equivalent programs in native code for the software switch. We demonstrated that, with appropriate compiler optimizations, this drastic reduction in complexity incurs only a small performance overhead compared to the native software switch implementation.

Our prototype demonstrates the feasibility of a protocol-independent software switch using P4 as the programming language and OVS as the target switch. Moreover, our techniques for software switch protocol independence and for compiling a domain-specific packet-processing language to an efficient low-level implementation should generalize to other languages and targets. One way to achieve language- and target-independence would be to first compile the domain-specific languages to a protocol-independent high-level intermediate representation (HLIR) such as protocol-oblivious forwarding [68] or NetASM [65], then apply the techniques and optimizations from PISCES to the HLIR.

Another future enhancement for PISCES is to enable custom parse, match, and action code to be dynamically loaded into a running protocol-independent switch. PISCES currently requires recompilation of the switch source code every time the programmer changes the P4 specification. In certain instances, such as adding new features and protocols to running production switches or temporarily altering protocol behavior to add visibility or defend against an attack, dynamically loading code in a running switch would be valuable. We expect future programmable protocol-independent software switches to sup-

port dynamically loading new or modified packet-processing code. Finally, PISCES does not implement P4 features that maintain state across packets (*i.e.*, counters, meters, or registers), which would require extending and generalizing the Open vSwitch caching model to achieve acceptable performance.

It is too early to see the effects of PISCES on protocol development, but the resulting code simplicity should make it easier to deploy, implement, and maintain custom software switches. In particular, protocol designers can maintain their custom software switch implementations in terms of a high-level domain-specific language like P4 without needing to track the evolution of the (larger and more complex) underlying software switch codebase. The ability to develop proprietary customizations without having to modify (and track) the source code for a software switch such as OVS might also be a selling point for protocol designers. We intend to study and characterize these effects as we release PISCES and interact with the protocol designers who use it.

Acknowledgments

We thank our shepherd Jeff Mogul, William Tu, and the anonymous SIGCOMM reviewers for their valuable feedback that helped improve the quality of this paper. We also thank Chaitanya Kodeboyina, Mihai Budiu, Ramkumar Krishnamoorthy, Antonin Bas, Abhinav Narain, and Bilal Anwer for their invaluable support at various stages of this project. This research was supported by Open Networking Research Center (ONRC), The Stanford Platform Lab, National Science Foundation (NSF) Awards CNS-1531281, CNS-1162112, and a generous gift from Intel.

References

- [1] P4 program for OVS, June 2015. <https://github.com/blp/ovs-reviews/blob/p4-workshop/tests/ovs.p4>.
- [2] P4-vSwitch. <https://github.com/P4-vSwitch>, 2016.
- [3] A. V. Aho, R. Sethi, and J. D. Ullman. *Compilers: Principles, Techniques, and Tools*. Addison-Wesley Longman Publishing Co., Inc., 1986.
- [4] M. Alizadeh, T. Edsall, S. Dharmapurikar, R. Vaidyanathan, K. Chu, A. Fingerhut, V. T. Lam, F. Matus, R. Pan, N. Yadav, and G. Varghese. CONGA: Distributed Congestion-aware Load Balancing for Datacenters. In *ACM SIGCOMM*, pages 503–514, 2014.
- [5] D. G. Andersen, H. Balakrishnan, N. Feamster, T. Koponen, D. Moon, and S. Shenker. Accountable Internet Protocol (AIP). In *ACM SIGCOMM*, pages 339–350, 2008.
- [6] M. B. Anwer, M. Motiwala, M. b. Tariq, and N. Feamster. SwitchBlade: A Platform for Rapid Deployment of Network Protocols on Programmable Hardware. In *ACM SIGCOMM*, pages 183–194, 2010.
- [7] G. Back. DataScript: A Specification and Scripting Language for Binary Data. In *ACM SIGPLAN/SIGSOFT*, pages 66–77. Springer-Verlag, 2002.

- [8] W. Bai, L. Chen, K. Chen, D. Han, C. Tian, and H. Wang. Information-agnostic Flow Scheduling for Commodity Data Centers. In *12th USENIX Conference on Networked Systems Design and Implementation (NSDI)*, pages 455–468, 2015.
- [9] A. Bianco, R. Birke, L. Giraudo, and M. Palacin. OpenFlow Switching: Data Plane Performance. In *IEEE International Conference on Communications (ICC)*, pages 1–5, 2010.
- [10] 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. *ACM SIGCOMM Computer Communication Review (CCR)*, 44(3):87–95, July 2014.
- [11] 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 *ACM SIGCOMM*, pages 99–110, 2013.
- [12] G. Brebner. Programmable Hardware for Software Defined Networks. In *IEEE European Conference on Optical Communication (ECOC)*, pages 1–3, 2015.
- [13] D. Coleman, D. Ash, B. Lowther, and P. Oman. Using Metrics to Evaluate Software System Maintainability. *IEEE Computer*, 27(8):44–49, 1994.
- [14] J. Corbet. BPF: The Universal In-kernel Virtual Machine. *Linux Weekly News, Eklektix Inc*, 2014.
- [15] J. Corbet. Extending BPF. *Linux Weekly News, Eklektix Inc*, 2014.
- [16] B. Davie and J. Gross. A Stateless Transport Tunneling Protocol for Network Virtualization (STT). Internet-Draft draft-davie-stt-08, Internet Engineering Task Force, Apr. 2016. Work in Progress.
- [17] M. Dillon and T. Winters. Network Functions Virtualization in Home Networks. Technical report, Open Networking Foundation, 2015. <https://www.opennetworking.org/images/stories/downloads/sdn-resources/IEEE-papers/network-func-virt-in-home-networks.pdf>.
- [18] M. Dobrescu, N. Egi, K. Argyraki, B.-G. Chun, K. Fall, G. Iannaccone, A. Knies, M. Manesh, and S. Ratnasamy. RouteBricks: Exploiting Parallelism to Scale Software Routers. In *ACM SIGOPS 22nd Symposium on Operating Systems Principles (SOSP)*, pages 15–28, 2009.
- [19] N. Dukkupati, G. Gibb, N. McKeown, and J. Zhu. Building a RCP (Rate Control Protocol) Test Network. In *15th IEEE Symposium on High-Performance Interconnects (HOTI)*, pages 91–98, 2007.
- [20] P. Emmerich, S. Gallenmüller, D. Raumer, F. Wohlfart, and G. Carle. MoonGen: A Scriptable High-Speed Packet Generator. In *ACM Internet Measurement Conference (IMC)*, pages 275–287, 2015.
- [21] P. Emmerich, D. Raumer, F. Wohlfart, and G. Carle. Performance Characteristics of Virtual Switching. In *IEEE International Conference on Cloud Networking (CloudNet)*, pages 120–125, 2014.
- [22] D. Farinacci, S. P. Hanks, D. Meyer, and P. S. Traina. Generic Routing Encapsulation (GRE). RFC 2784, Mar. 2000.
- [23] K. Fisher and R. Gruber. PADS: A Domain-specific Language for Processing Ad Hoc Data. In *ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, pages 295–304, 2005.
- [24] N. Foster, R. Harrison, M. J. Freedman, C. Monsanto, J. Rexford, A. Story, and D. Walker. Frenetic: A Network Programming Language. In *16th ACM SIGPLAN International Conference on Functional Programming (ICFP)*, pages 279–291, 2011.
- [25] T. M. Gil and M. Poletto. MULTOPS: A Data-structure for Bandwidth Attack Detection. In *10th Conference on USENIX Security Symposium*, 2001.
- [26] A. Greenberg, J. R. Hamilton, N. Jain, S. Kandula, C. Kim, P. Lahiri, D. A. Maltz, P. Patel, and S. Sengupta. VL2: A Scalable and Flexible Data Center Network. In *ACM SIGCOMM*, pages 51–62, 2009.
- [27] J. Gross. Tunnel: Add support for matching on OAM packets. Git commit 94872594b79d in [53], May 2014.
- [28] J. Gross. Tunneling: Allow matching and setting tunnel ‘OAM’ flag. Git commit b666962be3b2 in [53], July 2015.
- [29] J. Gross and I. Ganga. Geneve: Generic Network Virtualization Encapsulation. Internet-Draft draft-ietf-nvo3-geneve-01, Internet Engineering Task Force, Jan. 2016. Work in Progress.
- [30] S. Han, K. Jang, A. Panda, S. Palkar, D. Han, and S. Ratnasamy. SoftNIC: A Software NIC to Augment Hardware. Technical Report UCB/EECS-2015-155, EECS Department, University of California, Berkeley, May 2015.
- [31] S. Han, K. Jang, K. Park, and S. Moon. PacketShader: A GPU-accelerated Software Router. In *ACM SIGCOMM*, pages 195–206, 2010.
- [32] N. Heirbaut and T. Van Der Storm. Two implementation techniques for domain specific languages compared: OMeta/JS vs. JavaScript. *Master’s thesis, Universiteit van Amsterdam*, 2009.
- [33] D. Hiebert. Ctags User Commands Version 5.8-1. *Exuberant Ctags*.
- [34] Intel. DPDK: Data Plane Development Kit. <http://dpdk.org>, 2013.
- [35] Intel. DPDK: Data Plane Development Kit - Programmer’s Guide, 2013. http://dpdk.org/doc/guides/program_guide/index.html.
- [36] L. Jose, L. Yan, G. Varghese, and N. McKeown. Compiling Packet Programs to Reconfigurable Switches. In *12th USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, pages 103–115, 2015.
- [37] S. H. Kan. *Metrics and Models in Software Quality Engineering*. Addison-Wesley Longman Publishing Co., Inc., 2nd edition, 2002.
- [38] C. Kaner et al. Software engineering metrics: What do they measure and how do we know? In *IEEE METRICS*. Citeseer, 2004.
- [39] D. Katabi, M. Handley, and C. Rohrs. Congestion Control for High Bandwidth-delay Product Networks. In *ACM SIGCOMM*, pages 89–102, 2002.
- [40] N. Katta, M. Hira, C. Kim, A. Sivaraman, and J. Rexford. HULA: Scalable Load Balancing Using Programmable Data Planes. In *2nd ACM SIGCOMM Symposium on Software Defined Networking Research (SOSR)*, 2016.

- [41] C. Kim. Programming the Network Dataplane in P4, 2016. http://netseminar.stanford.edu/03_31_16.html.
- [42] C. Kim, P. Bhide, E. Doe, H. Holbrook, A. Ghanwani, D. Daly, M. Hira, and B. Davie. In-band Network Telemetry (INT), 2016. <http://p4.org/wp-content/uploads/fixed/INT/INT-current-spec.pdf>.
- [43] C. Kim, A. Sivaraman, N. Katta, A. Bas, A. Dixit, and L. J. Wobker. In-band Network Telemetry via Programmable Dataplanes. In *ACM SIGCOMM*, 2015. Demo Session.
- [44] E. Kohler, R. Morris, B. Chen, J. Jannotti, and M. F. Kaashoek. The Click Modular Router. *ACM Transaction on Computer Systems (TOCS)*, 18(3):263–297, Aug. 2000.
- [45] C. Lameter. NUMA (Non-Uniform Memory Access): An Overview. *ACM Queue*, 11(7):40, 2013.
- [46] Linux Kernel Archives. <http://kernel.org>, 1997.
- [47] M. Mahalingam, T. Sridhar, M. Bursell, L. Kreeger, C. Wright, K. Duda, P. Agarwal, and D. Dutt. Virtual eXtensible Local Area Network (VXLAN): A Framework for Overlaying Virtualized Layer 2 Networks over Layer 3 Networks. RFC 7348, Oct. 2015.
- [48] P. J. McCann and S. Chandra. Packet Types: Abstract Specification of Network Protocol Messages. In *ACM SIGCOMM*, pages 321–333, 2000.
- [49] S. McCanne and V. Jacobson. The BSD Packet Filter: A New Architecture for User-level Packet Capture. In *USENIX*, pages 2–2, 1993.
- [50] N. McKeown, T. Anderson, H. Balakrishnan, G. Parulkar, L. Peterson, J. Rexford, S. Shenker, and J. Turner. OpenFlow: Enabling Innovation in Campus Networks. *ACM SIGCOMM Computer Communication Review (CCR)*, 38(2):69–74, Mar. 2008.
- [51] Network Working Group. RFC 1624: Computation of the Internet Checksum via Incremental Update, May 1994.
- [52] P. Oman and J. Hagemeister. Metrics for assessing a software system’s maintainability. In *Conference on Software Maintenance*, pages 337–344, 1992.
- [53] Open vSwitch. <https://github.com/openvswitch/ovs>, October 2015.
- [54] R. Pagh and F. F. Rodler. Cuckoo hashing. *Elsevier Journal of Algorithms*, 51(2):122–144, 2004.
- [55] I. Pepelnjak. Packet Forwarding in Amazon VPC, December 2013. <http://blog.ipSPACE.net/2013/12/packet-forwarding-in-amazon-vpc.html>.
- [56] B. Pfaff. P4 Parsing in Open vSwitch, June 2015. P4 Workshop, <http://p4workshop2015.sched.org/event/3ZQF>.
- [57] 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 *12th USENIX Conference on Networked Systems Design and Implementation (NSDI)*, pages 117–130, 2015.
- [58] S. Previdi et al. *SPRING Problem Statement and Requirements*. IETF, June 2015. <https://datatracker.ietf.org/doc/draft-ietf-spring-problem-statement>.
- [59] Proxmox Virtual Environment. <https://www.proxmox.com/en/proxmox-ve>.
- [60] P. Quinn and U. Elzur. Network Service Header. Internet-Draft draft-ietf-sfc-nsh-04, Internet Engineering Task Force, Mar. 2016. Work in Progress.
- [61] J. Rajahalme. TCP flags matching support. Git commit dc235f7bcbff in [53], October 2013.
- [62] J. Reich, C. Monsanto, N. Foster, J. Rexford, and D. Walker. Modular SDN Programming with Pyretic. *USENIX ;login:*, 38(5):128–134, 2013.
- [63] T. Rinta-Aho, M. Karlstedt, and M. P. Desai. The Click2NetFPGA Toolchain. In *USENIX Annual Technical Conference (ATC)*, pages 7–7, 2012.
- [64] L. Rizzo. Netmap: A Novel Framework for Fast Packet I/O. In *USENIX Annual Technical Conference (ATC)*, pages 101–112, June 2012.
- [65] M. Shahbaz and N. Feamster. The Case for an Intermediate Representation for Programmable Data Planes. In *1st ACM SIGCOMM Symposium on Software Defined Networking Research (SOSR)*, pages 31–36, 2015.
- [66] N. Shelly, E. J. Jackson, T. Koponen, N. McKeown, and J. Rajahalme. Flow Caching for High Entropy Packet Fields. In *34th ACM SIGCOMM Workshop on Hot Topics in Software Defined Networking (HotSDN)*, pages 151–156, 2014.
- [67] M. Smith and L. Kreeger. VXLAN Group Policy Option. Internet-Draft draft-smith-vxlan-group-policy-02, Internet Engineering Task Force, Apr. 2016. Work in Progress.
- [68] H. Song. Protocol-oblivious Forwarding: Unleash the Power of SDN Through a Future-proof Forwarding Plane. In *2nd ACM SIGCOMM Workshop on Hot Topics in Software Defined Networking (HotSDN)*, pages 127–132, 2013.
- [69] V. Srinivasan, S. Suri, and G. Varghese. Packet Classification Using Tuple Space Search. In *ACM SIGCOMM*, pages 135–146, 1999.
- [70] J. Stringer. datapath: Allow matching on conntrack label. Git commit 038e34abaa31 in [53], December 2012.
- [71] J. Stringer. Add connection tracking label support. Git commit 9daf23484fb1 in [53], October 2013.
- [72] M. Varvello, R. Laufer, F. Zhang, and T. Lakshman. Multi-Layer Packet Classification with Graphics Processing Units. In *10th ACM International Conference on Emerging Networking Experiments and Technologies (CoNEXT)*, pages 109–120, 2014.
- [73] Y.-S. Wang and P. Garg. NVGRE: Network Virtualization Using Generic Routing Encapsulation. RFC 7637, Oct. 2015.
- [74] D. Zhou, B. Fan, H. Lim, M. Kaminsky, and D. G. Andersen. Scalable, High Performance Ethernet Forwarding with CuckooSwitch. In *9th ACM Conference on Emerging Networking Experiments and Technologies (CoNEXT)*, pages 97–108, 2013.