

# A Compiler and Run-time System for Network Programming Languages

Christopher Monsanto  
Princeton University

Nate Foster  
Cornell University

Rob Harrison \*  
United States  
Military Academy

David Walker  
Princeton University

## Abstract

Software-defined networks (SDNs) are a new implementation architecture in which a controller machine manages a distributed collection of switches, by instructing them to install or uninstall packet-forwarding rules and report traffic statistics. The recently formed Open Networking Consortium, whose members include Google, Facebook, Microsoft, Verizon, and others, hopes to use this architecture to transform the way enterprise and data center networks are implemented. But to do this, they need novel programming languages to help them craft network-wide algorithms for routing, energy-efficient network management, dynamic access control, traffic monitoring, and other applications.

In this paper, we define a high-level language, called NCore, for expressing packet-forwarding policies and traffic-statistics queries. The language is designed to be simple, expressive, and compositional. We define a formal semantics for NCore and show how to compile it to a distributed switch-controller architecture. To ensure that a majority of packets are processed efficiently on switches, we develop a new compilation technique called *reactive specialization* that generalizes, improves on, and automates the simple (but inefficient) manual techniques commonly used to program SDNs. Reactive specialization and the other compilation techniques we develop are highly generic, assuming only that the packet-matching capabilities available on switches satisfy some basic algebraic laws. This generality makes our technology applicable to all current switches we are aware of, including switches that implement the popular OpenFlow protocol.

Overall, this paper delivers a design for a high-level network programming language; a novel, general-purpose compilation algorithm based on reactive specialization; a run-time system based on a SDN architecture; the first formal semantics and proofs of correctness in this domain; and an implementation and evaluation that demonstrates the benefits over the current state-of-the-art.

## 1. Introduction

A network is a collection of connected devices that route traffic from one place to another. Networks are pervasive: they connect students and faculty on university campuses, they send packets between a variety of mobile devices in modern households, they route search requests and shopping orders through data centers, they tunnel between corporate networks in San Francisco and Helsinki, and they connect the steering wheel to the drive train in your car. Naturally, these networks have different purposes, properties, and requirements. To service these requirements,

companies like Cisco, Juniper, and others manufacture a variety of devices including routers (which forward packets based on IP addresses), switches (which forward packets based on MAC addresses), NAT boxes (which translate addresses within a network), firewalls (which squelch forbidden or unwanted traffic), and load balancers (which distribute work among servers), to name a few.

While each of these devices behaves differently, internally they are all built on top of a *data plane* that buffers, forwards, drops, tags, rate limits, and collects statistics about packets at high speed. More complicated devices like routers also have a *control plane* that run algorithms for tracking the topology of the network and computing routes through it. Using the statistics gathered from the data plane and the results computed using the device's specialized algorithms, the control plane installs or uninstalls new forwarding rules in the data plane. The data plane is built out of fast, special-purpose hardware, capable of forwarding packets at the rate at which they arrive, while the control plane is typically implemented in software.

Remarkably, however, traditional networks appear to be on the verge of a major upheaval. On March 11th, 2011, six companies, Deutsche Telekom, Facebook, Google, Microsoft, Verizon, and Yahoo!, owners of some of the largest networks in the world, announced the formation of the Open Networking Foundation [20]. The foundation's proposal is extraordinarily simple: *eliminate the control plane from network devices*. Rather than baking specific control software into each device, support a standard protocol that allows a separate, general-purpose machine called a *controller* to program and query the data planes of many cooperating devices. By moving the control plane from special-purpose devices onto stock machines, companies like Facebook and Google will be able to buy cheap, commodity switches, and write controller programs to customize and optimize their networks however they choose.

Networks built on this new architecture, which arose from earlier work on Ethane [4] and 4D [10], are now commonly referred to as *Software-Defined Networks* (SDNs). Already, a number of commercial switch vendors support OpenFlow [17], a concrete realization of the switch-controller protocol required for implementing SDNs, and many researchers have used OpenFlow to develop new network-wide algorithms for server load-balancing, data center routing, energy-efficient network management, virtualization, fine-grained access control, traffic monitoring, fault tolerance, denial of service detection, and host mobility [8, 12–14, 19, 27].

Now the obvious question is: Why should programming language researchers, and the POPL community in particular, care about these developments? The answer is clear: Some of our most important infrastructure—our networks—will soon be running *an entirely new kind of program*, and using our experience, principles, tools and algorithms, our community has a unique opportunity to define the languages these programs will be written in and the in-

\*The views expressed in this paper are those of the authors and do not reflect the official policy or position of the United States Military Academy, the Department of the Army, the Department of Defense, or the U.S. Government.

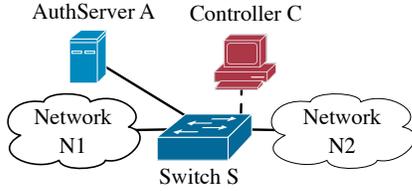


Figure 1. Example topology.

frastructure used to implement them. We can have major impact, making future networks more secure, reliable, and efficient.

As a step toward carrying out this agenda, we propose a high-level programming language called NCore, the *Network Core Programming Language*, for expressing packet-forwarding policies and statistics queries. NCore has an intuitive syntax based on familiar set-theoretic operations, which allows programmers to construct (and reason about!) rich policies in a natural way. NCore’s primitives include packet predicates, which are relatively easy to compile to conventional switch hardware, as well as black-box functions, which may depend on the state of the controller. Such functions make it possible to describe complicated, dynamic policies such as authentication and load balancing. Although they are not supported natively on any standard switches, black-box functions can be implemented in SDNs with help from the controller.

The NCore compiler analyzes programs and automatically divides them into two pieces: one that runs on the switches and another that runs on the controller. The main goal of the compiler is to arrange for as much packet processing as possible to occur on the switches because packets redirected to the controller are slowed down by several orders of magnitude. However, as suggested above, not all NCore program parts can be compiled into static packet-processing rules for switches—some, like authentication and load balancing, are fundamentally dynamic. To cope with these programs, we propose a new technique, called *reactive specialization*, that dynamically refines the NCore program using the traffic seen at run time. The additional information obtained from this refinement allows the run-time system to generate and install new, specialized packet-processing rules on switches. Hence, over time, more and more rules are added to the switches and less and less traffic is diverted to the controller.

Our strategy is inspired by the idiom commonly used for SDN applications today [12, 13, 28], in which an event-driven program manually installs a rule to handle future traffic every time a packet is diverted to the controller. However, we improve on it by (1) automating the process of dynamically unfolding packet-processing rules on to switches instead of requiring that programmers craft tricky, low-level, event-based programs manually, (2) synthesizing efficient forwarding rules that exploit the capabilities of modern switches (e.g., wildcard rules implemented using TCAMs), and (3) providing high-level abstractions that obviate the need for programmers to deal with the low-level details of individual switches.

To summarize: the central contribution of this paper is a framework for implementing a canonical, high-level network programming language correctly and efficiently. More specifically:

- We define a denotational semantics for NCore (Section 3) and model the interaction between the NCore run-time system and the network in a process calculus style (Section 4). Despite the huge industrial and academic momentum surrounding SDNs such as OpenFlow, no one has formalized the low-level details of how a controller system interacts with switches yet.

- We develop novel algorithms for compiling network programs, including *algebraic classifier generation* and *reactive specialization* (Section 5).
- We prove key correctness theorems concerning our compiler and run-time system (Section 6), establishing simulation relations between our low-level, distributed implementation strategy and our high-level NCore semantics. We also prove important *progress* and *quiescence* theorems, which show that our implementation successfully moves computation off of the controller and onto switches.
- We describe a prototype implementation and an evaluation on some simple benchmarks demonstrating the practical utility of our framework (Section 7).

## 2. SDN and NCore Overview

This section presents additional background on SDN and NCore, using a running example to illustrate the main ideas. For concreteness, we focus on the OpenFlow SDN architecture [21] (eliding and taking liberties with some of its inessential details). Note however, that our compiler is designed to be general, and does not assume the specifics of the current OpenFlow platform.

**OpenFlow overview.** OpenFlow is based on a two-tiered architecture in which a controller manages a collection of subordinate switches. Figure 1 depicts a simple topology with a controller  $C$  managing a single switch  $S$ . Packets may either be processed on the switches or on the controller, but processing a packet on the controller increases its latency by several orders of magnitude. Hence, to ensure good performance, the controller typically installs a *classifier* on each switch comprising a list of *rules* that can be used to process packets directly on the switch, without having to send them to the controller.

Each forwarding rule has a *pattern* that identifies a set of packets, a list of *actions* that specifies how packets matching the pattern should be processed, a *counters* that keep track of the number and size of all packets processed using the rule, and an integer *priority*. When a packet arrives at a switch, it is processed in three steps: First, the switch selects a rule whose pattern matches the packet. If it has no matching rules, then it drops the packet, and if it has multiple matching rules, then it picks the one with the highest priority. Second, the switch updates the counters associated with the rule. Finally, the switch applies each of the actions listed in the rule to the packet. OpenFlow supports two kinds of actions:  $\text{modify}(h, n)$ , which sets the value of header field  $h$  to  $n$ , and  $\text{output}(l)$ , which forwards the packet to an adjacent network location  $l$ . A *location*  $l$  is either the name of a switch, network, or host, or the special location *controller*.<sup>1</sup>

**NCore Example.** Suppose that we want to build a security application that implements in-network authentication for the network in Figure 1. The network  $N_1$  contains a collection of internal hosts,  $N_2$  represents the upstream connection to the Internet,  $A$  is the server that handles authentication for hosts in  $N_1$ , and all three elements are connected to each other by the switch  $S$ .

Informally, we want the network to perform routing and access control according to the following policy: Forward packets from unauthenticated hosts in  $N_1$  to  $A$ , from authenticated hosts in  $N_1$  to their intended destination in  $A$  or in  $N_2$ , and from  $A$  and  $N_2$  back to  $N_1$  (although not from  $N_2$  to  $A$ ). We will assume that once a host is authenticated, it is authenticated for all time.

This policy can be described succinctly in NCore as follows:

<sup>1</sup> On real OpenFlow switches, locations are actually integers corresponding to physical ports on the switch; in this paper we model them symbolically.

```

((InPort <: Network 1) /\ inspect auth_bb
 --> [Forward (Network 2)])
\ / ((InPort <: Network 1) /\ neg (inspect auth_bb)
 --> [Forward (AuthServer A)])
\ / ((InPort <: AuthServer A) \ / (InPort <: Network 3)
 --> [Forward (Network 1)])

```

Intuitively, this policy can be read as a function from triples consisting of the controller state, a switch, and a packet, to sets of actions. Formally, it is structured as the disjunction of three clauses each consisting of a *predicate* describing a set of packets (at a particular switch and controller state) and a set of *actions*, written concretely as a list. The first states that packets arriving from the internal network (`InPort <: Network 1`) should be forwarded to the external network (`Forward (Network 2)`) if `auth_bb` holds in the current state. The `inspect` keyword wraps a user-defined predicate, in this instance, `auth_bb`, which tests whether the host listed as the source in the packet header is authenticated. Similarly, the second clause states that packets from the internal network should be forwarded to the authentication server if `auth_bb` does not hold. The final clause states that packets from the authentication server and external network should be forwarded to the internal network.

Besides specifying the forwarding behavior of the network, administrators often need to monitor it for security purposes, or for making decisions about load balancing, energy efficiency, and billing. To keep things simple, in this example we will monitor the traffic from the authentication server, aggregating the results by destination host, as well as the non-web traffic (*i.e.*, not port 80) from internal hosts to the external network, aggregated by source host. We can define a query in NCore that monitors these traffic statistics, using *keys* to keep track of each result. For the first, we introduce a key `auth` and assign it the indexed type `Key[{dstmac}]`, which indicates that only the `dstmac` field from each packet should be included in the result. Likewise, for the second query, we call the key `traf` and associate it with the type `Key[{srcmac}]`. Using these keys, the query itself can be defined as follows:

```

(auth <-- (InPort <: AuthServer A))
\ / (traf <-- (InPort <: Network 1) /\ neg (DstPort <: 80))

```

The overall program consists of the forwarding policy and the query.

Unfortunately, although these programs are easy to write, they are far from easy to implement—at least, not in an efficient way. Compiling these programs to an efficient distributed implementation on top of an OpenFlow-like architecture requires solving two challenging problems.

**Challenge 1: Dynamic program partitioning.** OpenFlow’s architecture is asymmetric—the controller and the switches have vastly different capabilities. To implement programs efficiently, the compiler must maximize the amount of processing done on the fast (but computationally limited) switches and minimize the amount done on the slow (but omnipotent) controller. In general, it cannot simply push the entire computation onto the switches because the switch hardware is fundamentally limited—it can match bit patterns, maintain counters, and execute modify/forward actions, but it cannot directly evaluate functions like `auth_bb` or collect statistics like the total amount of non-web traffic per host (without several billion rules). Hence, a compiler and run-time system for an SDN must solve the fundamental problem of determining what parts of the overall computation to execute where.

In general, as programs execute, they gain information that can be exploited to push work onto the switches. Thus, although the complete space of possible packets is huge, the actual set of packets that will appear on any given network is far smaller, and those packets are also highly correlated with each other—a packet from one host is likely to be followed by many more packets from the

same host. Thus, although it would be impractical for the controller to push 4.3 billion rules out to the switches in advance (one for each possible IP address) to implement a fine-grained query such as the one that counts per-host non-web traffic, it is feasible for the controller to react to new packets from a particular host and dynamically generate a new classifier that tabulates statistics for that host using the counters on the switch. Similarly, although the controller cannot evaluate an arbitrary function like `auth_bb` on every possible packet in advance, when a packet arrives at the controller, it can evaluate the function on that particular packet, and, if the result is known to be fixed (*i.e. invariant*) on some set of similar packets, then it can generate a new classifier that handles future packets in that set directly on the switch.

The idea of dynamically reacting to traffic seen in the network is not new—it is used in all SDN applications (*e.g.*, [12, 13, 28]). However, in current systems, the programs that dynamically react to packets are crafted manually by programmers. For example, in NOX [11], a popular OpenFlow controller platform, programs are written in an event-driven style and explicitly react to the arrival of packets at the controller by running imperative routines that install or uninstall individual rules on switches. In contrast, examining the NCore example presented above, the reader will notice that there are no explicit calls to install or uninstall rules; there is not even a separate notion of switch code and controller code. There is just one high-level, centralized, declarative specification of the policy and queries. Hence, the compiler and run-time system must figure how to distribute the computation across the network efficiently, incrementally, and adaptively over time.

**Solution 1: Reactive Specialization.** *Reactive specialization* is an implementation strategy for NCore that automates the construction of the kinds of reactive programs that SDN programmers currently write by hand. This frees programmers from having to worry about the low-level details of when to install or uninstall rules on switches and—by construction—prevents them from making mistakes in doing so.

When a packet arrives at the controller, the NCore run-time automatically refines the program, using the packet as a guide. For example, returning to the authentication program, if a packet `p` from source MAC `m` arrives at the controller from an authenticated host, then `auth_bb` holds, and, moreover, will continue to hold in the future. Using this information, the compiler transforms the `inspect auth_bb` predicate into a disjunction of two predicates:

```
(inspect auth_bb) \ / (SrcMac <: MAC m)
```

After applying this program transformation, the compiler has sufficient information to generate and install new rules on the switch that handle all future traffic involving `m`.

This program transformation is enabled by the observation that each host goes through two distinct phases: unauthenticated and authenticated. When the host is unauthenticated, the forwarding policy can only be determined by the controller because it might change (and the switches do not have the computational facilities necessary to recognize the change). But when it is authenticated, the effective forwarding policy for that host suddenly becomes *consistent* and *determinate*—no further state-based decisions are needed and all future packets can be forwarded directly by switches to their intended destination. Thus, the controller can safely install rules that do this.

**Challenge 2: Switch rule selection.** A second challenge that makes compilation difficult is the large gap between the languages used to express network policies and switch-level classifiers. With that in mind, consider the following policy:<sup>2</sup>

<sup>2</sup>We use standard conventions when writing IP prefixes in our code—*e.g.*, `10.0.0.0/8` denotes the network of IP addresses whose first octet is 10.

```
(SrcAddr <: IP "10.0.0.0/8") //
(SrcAddr <: IP "10.0.0.1" \\/ DstPort <: Port 80)
--> [Forward (Switch 1)]
```

It states that all packets from sources in the network identified by 10.0.0.0/8 should be forwarded to switch 1, except for packets coming from 10.0.0.1 or going to a destination on port 80 (the // operator denotes set-theoretic difference). Because switches cannot express the difference of two patterns in a single rule, this policy needs to be implemented using three rules installed in a particular order: one that drops packets from 10.0.0.1, another that drops all packets going to port 80, and a final rule that forwards all remaining packets from 10.0.0.0/8 to switch 1. The following classifier, which has the highest priority rule placed first, implements this policy:

```
SrcAddr 10.0.0.1      : []
DstPort 80           : []
SrcAddr 10.0.0.0/8   : [Forward 1]
```

Next consider a similar policy:

```
(SrcAddr <: IP "10.2.0.0/16") //
(SrcAddr <: IP "10.2.0.1" \\/ DstPort <: Port 22)
--> [Forward (Switch 2)]
```

We can generate a classifier for this policy in the same way:

```
SrcAddr 10.2.0.1      : []
DstPort 22           : []
SrcAddr 10.2.0.0/16   : [Forward 2]
```

Now suppose that we then want to generate a classifier that implements the union of the two policies. We cannot combine the classifiers in a simple way—*e.g.*, by concatenating or interleaving them—because the compiled rules interact with each other. For example, if we were to simply concatenate the two lists of rules, the rule that drops packets to port 80 would incorrectly shadow the forwarding rule for traffic from 10.2.0.0/16. Instead, we need to perform a much more complicated translation that produces the following classifier:

```
SrcAddr 10.2.0.1,      DstPort 80 : [],
SrcAddr 10.2.0.0/16, DstPort 80 : [Forward 2]
SrcAddr 10.2.0.1      : [Forward 1]
SrcAddr 10.2.0.0/16, DstPort 22 : [Forward 1]
SrcAddr 10.2.0.0/16   : [Forward 1,Forward 2]
SrcAddr 10.0.0.1      : []
SrcAddr 10.0.0.0/8,   DstPort 80 : []
SrcAddr 10.0.0.0/8    : [Forward 1]
```

Dealing with these complexities often leads SDN programmers to use *exact-match* rules—*i.e.*, rules that specify a value for *every single header field*, without using larger rules whose patterns contain wildcards like 10.0.0.0/8 or leave some header fields unconstrained. As interactions between exact-match rules are easy to reason about, composing policies is straightforward. But, exact-match rules, which match narrow bands of traffic are far less efficient than wildcard rules. The results of our experiments, presented in Section 7, highlight the magnitude of the difference.

**Solution 2: Algebraic classifier generation.** Our solution to the second challenge is a general algorithm that automates the task of generating classifiers. Given an NCore program, this algorithm constructs a low-level classifier that implements the program as fully and efficiently as possible. The classifier is efficient because the algorithm automatically selects and overlays wildcard rules whenever possible. However, because of the presence of black box functions and unimplementable queries, some of the rules will send packets to the controller for further processing instead of sending them directly to their destinations.

Importantly, our classifier generation algorithm interacts harmoniously with reactive specialization. This is critical because at run

time, the output of the latter is supplied as the input to the former. Thus, the classifier generation algorithm must be powerful enough to extract new information from the program produced by reactive specialization and generate new classifiers that send fewer packets back to the controller. The theoretical analysis presented in Section 6, especially the *progress* and *quiescence* theorems, demonstrates that our algorithms interoperate harmoniously.

### 3. A Core Calculus for Network Programming

This section defines the syntax and semantics of NCore, a core calculus for high-level network programming. The calculus has three major components: *predicates*, which describe sets of packets; *policies*, which associate high-level *actions* with the predicates that trigger them; and *queries*, which describe various statistics about the packets flowing through the network. Figure 2 presents the syntax of these constructs as well as various network values such as headers and packets.

**Notation.** Throughout this paper, whenever we define syntax, as in the grammar for packets, we will use the grammar non-terminal ( $p$ ) as a metavariable ranging over the objects being defined, the capitalized version of the non-terminal ( $P$ ) as a metavariable ranging over sets of such objects, the overlined and capitalized version ( $\overline{P}$ ) for multisets of objects, and vector notation ( $\vec{p}$ ) for sequences of objects. To convert a given set ( $P$ ) to the corresponding multiset, we simply place a line over the metavariable ( $\overline{P}$ ). We will use a name written in a sans-serif font to the left of the grammar (Packet) as the type of the objects and the name enclosed in curly braces ( $\{\text{Packet}\}$ ) as the type of sets of objects. Many of our types  $T$  are indexed by a set of headers  $H$ . We write such types  $T[H]$ . We write finite maps using the notation  $\{x_1 \mapsto y_1, \dots, x_n \mapsto y_n\}$  and we lookup elements of a map  $m$  using function application  $m(x_i)$ . We write multisets as  $\{\{x_1, \dots, x_n\}\}$  and take the union of two multisets  $M_1$  and  $M_2$  using the notation  $M_1 \uplus M_2$ . We use the usual operations  $\cup$ ,  $\cap$ ,  $\neg$  and  $\setminus$  for the operations of union, intersection, complement and difference on sets. Typically, we give definitions for intersection and complement and leave union ( $S_1 \cup S_2 = \neg(\neg S_1 \cap \neg S_2)$ ) and difference ( $S_1 \setminus S_2 = S_1 \cap \neg S_2$ ) as derived forms. We also overload  $\neg$  and use it to negate booleans; its meaning will be clear from context.

**Network Values.** For simplicity, we only model a single kind of location in our theory, switches  $s$ . Packets  $p$  are the basic values processed by programs. A *complete packet* has  $\ell$  header fields,  $h_1$  to  $h_\ell$ . We represent a packet as a finite map from these headers to bitstrings  $b$  and write  $p(h)$  for the bitstring associated with the header  $h$  in  $p$ . We assume all fields have fixed, finite length and therefore the set of complete packets is finite. A *subpacket*  $sp$  is like a packet, but is only defined on a subset  $H$  of the possible headers. We write  $p[H]$  for the subpacket obtained by restricting  $p$  to the headers in  $H$  and  $\text{Packet}[H]$  for the type of subpackets defined on  $H$ . A transmission  $t$  represents a packet in flight and is represented using the constructor  $\mathbb{T}$ .

**Predicates.** Informally, a predicate defines the conditions under which actions or queries are executed. Formally, a predicate  $e$  defines a set of *snapshots*  $x$  consisting of a controller state  $\Sigma$ , switch  $s$ , and packet  $p$ . The state and switch components are essential for modeling predicates defined in terms of black-box functions, such as the `auth_bb` function from Section 2. We say that a snapshot  $x$  *matches* a predicate  $e$  when it belongs to the denotation of  $e$ . We sometimes say that a packet  $p$  *matches*  $e$ , leaving the state and the switch implicit because they are uninteresting.

Figure 3 defines the semantics of predicates. Basic predicates have the form  $h : \vec{w}$ , where  $\vec{w}$  is a vector of wildcard bits. A bit  $b$  *matches* a wildcard bit  $w$ , written  $b \sqsubseteq w$ , if  $w$  is equal to  $b$  or

Network Values	
Switch	$s$
Header	$h$
Bit	$b ::= 1 \mid 0$
Packet	$p ::= \{h_1 \mapsto \bar{b}_1, \dots, h_\ell \mapsto \bar{b}_\ell\}$
Packet[H]	$sp ::= \{h_1 \mapsto \bar{b}_1, \dots, h_k \mapsto \bar{b}_k\}$
Transmission	$t ::= \mathbb{T}(s \mid p)$
Predicates	
Snapshot	$x ::= (\Sigma, s, p)$
WildcardBit	$w ::= 1 \mid 0 \mid ?$
Inspector	$i \in \text{Snapshot}[H] \rightarrow \text{Bool}$
Predicate	$e ::= h : \bar{w} \mid \text{switch } s \mid \text{inspect } i \mid e_1 \cap e_2 \mid \neg e$
Policies	
Action	$a ::= \text{forward } s \mid \text{modify } h \leftarrow \bar{b} \text{ in } A$
{Action}	$A ::= \{a_1, \dots, a_n\}$
Policy	$\tau ::= e \rightarrow A \mid \tau_1 \cap \tau_2 \mid \neg \tau$
Queries	
Key	$k \in \text{Key}[H]$
{Key}	$K ::= \{k_1, \dots, k_n\}$
PreState	$\sigma ::= \{(s_1, k_1) \mapsto SP_1, \dots, (s_n, k_n) \mapsto SP_n\}$
State	$\Sigma ::= \{(s_1, k_1) \mapsto \overline{SP}_1, \dots, (s_n, k_n) \mapsto \overline{SP}_n\}$
Query	$\phi ::= K \leftarrow e \mid \phi_1 \cap \phi_2 \mid \neg \phi$
Programs	
Program	$\pi ::= (\tau, \phi)$

Figure 2. NCore syntax.

? A vector of bits  $\bar{b}$  *matches* a vector of wildcard bits  $\bar{w}$ , written  $\bar{b} \sqsubseteq \bar{w}$ , if each bit in  $\bar{b}$  matches the corresponding wildcard bit in  $\bar{w}$ . A packet  $p$  matches a primitive predicate  $h : \bar{w}$  if the  $h$  header of  $p$  (that is,  $p(h)$ ) matches  $\bar{w}$ . Another basic predicate, *switch*  $s$ , denotes snapshots  $x$  containing  $s$ , with arbitrary state and packet components.

Black-box predicates, written *inspect*  $i$ , use an arbitrary function  $i$  to determine whether a snapshot  $x$  matches the predicate. The behavior of inspectors is constrained in two essential ways, which provides the compiler with sufficient information to generate switch-level rules. First, if  $i$  has type  $\text{Snapshot}[H] \rightarrow \text{Bool}$ , then the inspector may only examine the headers  $H$  of the packet contained in  $x$ . Second, each inspector has a (programmer supplied) *invariance oracle* that determines if the behavior of the inspector on a future snapshot with the same packet and switch will be the same as for  $x$ . More formally, an inspector  $i$  is *invariant on*  $(\Sigma, s, p)$ , written *invariant*  $((\Sigma, s, p), i)$ , if for all  $\Sigma'$ , we have  $i(\Sigma \uplus \Sigma', s, p) = i(\Sigma, s, p)$ , where  $\Sigma \uplus \Sigma'$  is multiset union applied pointwise to the underlying multisets contained within each state. As an example, the invariance oracle for the `auth_bb` inspector defined in Section 2 can be defined by the following function in Haskell notation:

```
\(state,s,p) ->
  is_auth state p || (dstaddr p) == authsrv_addr
```

This oracle promises that the `auth_bb` function will return the same answer on all packets coming from authorized hosts or going to the authentication server. Together, the inspector type and invari-

$\llbracket e \rrbracket \in \{\text{Snapshot}\}$	
$\llbracket h : \bar{w} \rrbracket$	$= \{(\Sigma, s, p) \mid p(h) \sqsubseteq \bar{w}\}$
$\llbracket \text{switch } s' \rrbracket$	$= \{(\Sigma, s, p) \mid s' = s\}$
$\llbracket \text{inspect } i \rrbracket$	$= \{(\Sigma, s, p) \mid i(\Sigma, s, p)\}$
$\llbracket e_1 \cap e_2 \rrbracket$	$= \llbracket e_1 \rrbracket \cap \llbracket e_2 \rrbracket$
$\llbracket \neg e \rrbracket$	$= \neg \llbracket e \rrbracket$
$\llbracket A \rrbracket \in \text{Packet} \rightarrow \{\text{Transmission}\}$	
$\llbracket \{a_1, \dots, a_n\} \rrbracket(p)$	$= T_1 \cup \dots \cup T_n$
where $T_i =$	$\begin{cases} \{\mathbb{T}(s \mid p)\} & \text{if } a_i = \text{forward } s \\ \llbracket A \rrbracket(p[h \mapsto \bar{b}]) & \text{if } a_i = \text{modify } h \leftarrow \bar{b} \text{ in } A \end{cases}$
$\llbracket \tau \rrbracket \in \text{Snapshot} \rightarrow \{\text{Transmission}\}$	
$\llbracket e \rightarrow A \rrbracket(\Sigma, s, p)$	$= \begin{cases} \llbracket A \rrbracket(p) & \text{if } (\Sigma, s, p) \in \llbracket e \rrbracket \\ \emptyset & \text{otherwise} \end{cases}$
$\llbracket \tau_1 \cap \tau_2 \rrbracket(\Sigma, s, p)$	$= \llbracket \tau_1 \rrbracket(\Sigma, s, p) \cap \llbracket \tau_2 \rrbracket(\Sigma, s, p)$
$\llbracket \neg \tau \rrbracket(\Sigma, s, p)$	$= \neg \llbracket \tau \rrbracket(\Sigma, s, p)$
$\llbracket K \rrbracket \in \text{Snapshot} \rightarrow \text{PreState}$	
$\llbracket \{k_1, \dots, k_n\} \rrbracket(\Sigma, s, p)$	$= \{(s, k_i) \mapsto \{p[H_i]\} \mid 1 \leq i \leq n\}$
where $k_i \in \text{Key}[H_i]$	
$\llbracket \phi \rrbracket \in \text{Snapshot} \rightarrow \text{PreState}$	
$\llbracket K \leftarrow e \rrbracket(x)$	$= \begin{cases} \llbracket K \rrbracket(x) & \text{if } x \in \llbracket e \rrbracket \\ \emptyset & \text{otherwise} \end{cases}$
$\llbracket \phi_1 \cap \phi_2 \rrbracket(x)$	$= \llbracket \phi_1 \rrbracket(x) \cap \llbracket \phi_2 \rrbracket(x)$
$\llbracket \neg \phi \rrbracket(x)$	$= \neg \llbracket \phi \rrbracket(x)$

Figure 3. NCore semantics.

ance oracle allow the compiler to generate effective switch-level rules even though the inspector function itself cannot be analyzed.

Complex predicates are built up from basic ones using the intersection and complement operators, which are interpreted set theoretically. Although this predicate language is minimal, it is straightforward to define additional, familiar operators—*e.g.*, `True` can be defined as  $h : ? \dots ?$  (with  $k$  question marks) for a specific  $k$ -width field  $h$ , `False` as  $\neg \text{True}$ , and  $e_1 \cup e_2$  as  $\neg(\neg e_1 \cap \neg e_2)$ .

**Policies.** Policies  $\tau$  specify how packets should be forwarded through the network. Basic policies, written  $e \rightarrow A$ , associate a set of actions  $A$  with a predicate  $e$ . If  $(\Sigma, s, p)$  matches  $e$ , where  $\Sigma$  is the current controller state, then this policy states that actions  $A$  should be applied to  $p$  at  $s$ . An action may either be a simple forwarding action *forward*  $s'$  that sends the packet to another switch  $s'$ , or a modification action (*modify*  $h \leftarrow \bar{b}$  in  $a$ ), which overwrites the value in header  $h$  with the vector  $\bar{b}$  and executes the actions  $A$  with respect to the modified packet. As with predicates, we build complex policies by combining simple policies using intersection and negation. Although negation is not useful on its own, as it can cause the machine to send every possible packet to every possible switch, we include it in the calculus for reasons of symmetry with predicates, and likewise can use it to encode useful operators on policies such as union and difference. Figure 3 defines the semantics of policies formally, as a function from snapshots  $x$  to sets of transmissions  $T$ .

Note that although the policy language is syntactically simple, it is remarkably expressive. In particular, inspectors, which are arbitrary functions that possibly depend upon state, are a powerful tool that can be used to express a wide range of policies. For example, it is easy to define policies that implement randomized and round-robin load balancing, fine-grained access control, and many standard routing protocols.

**Queries.** A query  $\phi$  selects a set of packets, projects out certain specified fields of those packets, yielding subpackets, and aggregates the subpackets under a specified set of keys. Basic queries ( $K \leftarrow e$ ) select packets that match  $e$  and associate them under each key in  $K$ . The projection is governed by the type of the key—if  $k$  belongs to the set  $K$  and  $k$  has type  $\text{Key}[H]$ , then  $p[H]$  is stored under  $k$ . In definitions, we write  $k \in \text{Key}[H]$  to indicate the type of  $k$  and  $K \in \text{KeySet}[H]$  where  $H$  is the union of the sets of headers  $H_i$  associated with each key  $k_i \in K$ . Figure 3 defines the semantics of queries as a function from snapshots to pre-states, which are maps from key-switch pairs to sets of subpackets. (States, used later to aggregate the results of queries, are maps from key-switch pairs to multisets of subpackets.) We interpret the operators  $\cap$  and  $\neg$  on pre-states point-wise.

**Programs.** An NCore program  $\pi$  consists of a pair of a policy  $\tau$  and query  $\phi$ . The semantics of a program on a given snapshot  $x$  is defined in the obvious way:  $\llbracket (\tau, \phi) \rrbracket (x) = (\llbracket \tau \rrbracket (x), \llbracket \phi \rrbracket (x))$ .

To understand how such a program executes over time, we define two abstract machines. Both machines forwards packets according to  $\tau$  and update the controller state using  $\phi$ , but they differ in how often the switches synchronize with the controller. The *synchronous machine* defines an idealized, synchronous implementation of an NCore program. It evaluates the query and policy on each packet immediately and instantaneously. Of course, it would be impractical to implement this machine in a real network because, in general, it would require sending every packet to the controller. The *asynchronous machine* defines a more realistic implementation. Like the first machine, it is *policy-compliant*—it always forwards packets according to the program policy—but it updates the state asynchronously rather than in lockstep with every packet processed. While the synchronous machine can be thought of as the best possible policy-compliant machine, the asynchronous machine can be thought of as the worst policy-compliant machine. Any reasonable implementation will sit between the two. In other words, a reasonable implementation should be policy compliant, but users should not expect perfect synchrony because the cost of implementing it would be prohibitive. In practice, synchronization typically happens at periodic, timed intervals (modulo the variance in the latency of communication) but for simplicity, we do not model time explicitly in our formal system.

Figure 4 defines both reference machines formally. The state of the synchronous machine includes the NCore program  $\pi$ , the state  $\Sigma$ , and a multiset  $\bar{T}$  of pending transmissions. At each step, it removes a transmission  $\bar{T}$ , processes it using the policy and query, and updates the machine state and transmissions with the new state and transmissions generated by the program. The state of the asynchronous machine includes the program  $\pi$ , state  $\Sigma$ , and two multisets of transmissions:  $\bar{T}_1$ , which represents transmissions waiting to be processed by the policy, and  $\bar{T}_2$ , which represents transmissions that have been processed by the policy but not yet by the query. The first inference rule for the second machine takes a transmission from  $\bar{T}_1$ , processes it using the policy, and places it in  $\bar{T}_2$ ; the second rule takes a transmission in  $\bar{T}_2$  and processes it using the query.

Before completing this section, we comment on one limitation of the language: the exclusive use of *permanent* invariance. There are policies that are invariant for a long time, and hence could have

$$\begin{array}{c}
 \boxed{M_{\text{sync}} \xrightarrow{o} M'_{\text{sync}}} \quad \text{SyncState } M_{\text{sync}} ::= (\pi, \Sigma, \bar{T}) \\
 \frac{\pi = (\tau, \phi) \quad \llbracket \tau \rrbracket (\Sigma, s, p) = T' \quad \llbracket \phi \rrbracket (\Sigma, s, p) = \sigma'}{(\pi, \Sigma, \bar{T} \uplus \{\mathbb{T}(s | p)\}) \xrightarrow{(s,p)} (\pi, \Sigma \uplus \bar{\sigma}', \bar{T} \uplus T')} \\
 \\
 \boxed{M_{\text{async}} \xrightarrow{o} M'_{\text{async}}} \quad \text{AsyncState } M_{\text{async}} ::= (\pi, \Sigma, \bar{T}_1, \bar{T}_2) \\
 \frac{\pi = (\tau, \phi) \quad \llbracket \tau \rrbracket (\Sigma, s, p) = T'}{(\pi, \Sigma, \bar{T}_1 \uplus \{\mathbb{T}(s | p)\}, \bar{T}_2) \xrightarrow{(s,p)} (\pi, \Sigma, \bar{T}_1 \uplus \bar{T}', \bar{T}_2 \uplus \{\mathbb{T}(s | p)\})} \\
 \\
 \frac{\pi = (\tau, \phi) \quad \llbracket \phi \rrbracket (\Sigma, s, p) = \sigma}{(\pi, \Sigma, \bar{T}_1, \bar{T}_2 \uplus \{\mathbb{T}(s | p)\}) \rightarrow (\pi, \Sigma \uplus \bar{\sigma}, \bar{T}_1, \bar{T}_2)}
 \end{array}$$

Figure 4. Reference machines.

rules installed on switches for that time, but are not permanently invariant. We believe our framework can be extended to handle such semi-permanent invariance properties, having the compiler uninstall rules at the end of the period, but a rigorous investigation of this topic is beyond the scope of this paper.

## 4. The Run-time System

The previous section defined the syntax and semantics of the high-level network programming language NCore. However, it did not discuss how to implement the semantics on a software-defined network—*i.e.*, on a distributed set of switches managed by a controller. In this section, we give an operational semantics to the NCore run-time system and the underlying network devices. This operational semantics explains the basic interactions between the controller and the switches. The next section will explain how to compile NCore programs and will use this low-level distributed system as the target of the compilation.

**Switch Classifiers.** Before we can present the run-time system, we need a concrete representation of the rules that switches use to process packets. A *classifier*  $\bar{r}$  is a sequence of rules  $r$  each containing a switch-level pattern  $z$  and action  $\gamma$ . We represent classifiers as sequences as opposed to sets primarily to model the features of modern switch hardware, including priorities (rules on the left have higher priority than rules on the right). Note that classifiers are more restricted than NCore policies—*e.g.*, individual rules cannot express negations or unions directly.

The *pattern* ( $z$ ) component of a rule recognizes a set of packets, and hence is similar to (but less general than) a predicate in NCore. We write  $p \sqsubseteq z$  when pattern  $z$  matches packet  $p$ . We assume that patterns form a lattice with elements higher in the lattice matching more packets and elements lower in the lattice matching fewer packets. We write  $z_1 \sqsubseteq z_2$  when  $z_1$  is lower than (or equal to) another element in the lattice, thereby matching a subset (or equal set) of packets. Operators  $\sqcap$  and  $\sqcup$  implement meets and joins in the pattern lattice. The top element ( $\top$ ) matches all packets; the bottom element ( $\perp$ ) matches none. In addition to  $\top$  and  $\perp$  and the pattern lattice contains a number of other elements. For example, OpenFlow switches support prefix pattern matching on source and destination IP addresses (*i.e.*, patterns of the form 0110\*) but only exact or wildcard matching on other most other headers. Other switches support alternative extended patterns, such as ranges with the form  $[n_1, n_2]$ . We do not specify the exact features available,

Pattern	$z ::= \top \mid \perp \mid \dots$	$\overline{M} \xrightarrow{o} \overline{M}'$			
Skeleton	$\vec{z} ::= (z_1, \dots, z_n)$		E-SWITCHHELP	E-SWITCHDROP	E-SWITCHPROCESS
SwitchAct	$\gamma ::= \Omega \mid \dots$		$\frac{\vec{r} \rightsquigarrow^p z : \Omega}{\mathbb{S}(s \vec{r} \vec{Z}), \mathbb{T}(s p) \rightarrow \mathbb{S}(s \vec{r} \vec{Z}), \mathbb{H}(s p)}$	$\frac{\vec{r} \rightsquigarrow^p}{\mathbb{S}(s \vec{r} \vec{Z}), \mathbb{T}(s p) \xrightarrow{s,p} \mathbb{S}(s \vec{r} \vec{Z})}$	$\frac{\vec{r} \rightsquigarrow^p z : \gamma \quad \gamma \neq \Omega}{\text{Action}(\gamma, p) = T'}$
Rule	$r ::= z : \gamma$				$\frac{\mathbb{S}(s \vec{r} \vec{Z}), \mathbb{T}(s p) \xrightarrow{s,p}}{\mathbb{S}(s \vec{r} \vec{Z} \uplus \{z\}), T'}$
Classifier	$\vec{r} ::= (r_1, \dots, r_n)$				
Help	$l ::= \mathbb{H}(s p)$	E-CONTROLLER	E-COLLECT		E-STEP
Molecule	$m ::= \mathbb{C}(\pi \Sigma) \mid \mathbb{S}(s \vec{r} \vec{Z}) \mid t \mid l$	$\frac{\llbracket \pi \rrbracket(\Sigma, s, p) = (T', \sigma')}{\text{Specialize}(\Sigma, s, p, \pi) = \vec{r}'}$	$\frac{\vec{r} \rightsquigarrow^p z : \gamma}{\llbracket \pi \rrbracket(\Sigma, s, p) = (T', \sigma')}$		$\frac{\overline{M} \xrightarrow{o} \overline{M}'}{\overline{M} \uplus \overline{M}'' \xrightarrow{o} \overline{M}' \uplus \overline{M}''}$
Observation	$o ::= \cdot \mid s, p$	$\frac{\mathbb{C}(\pi \Sigma), \mathbb{S}(s \vec{r} \vec{Z}), \mathbb{H}(s p) \xrightarrow{s,p}}{\mathbb{C}(\pi \Sigma \uplus \sigma'), \mathbb{S}(s \vec{r}', \vec{r} \vec{Z}), T'}$	$\frac{\mathbb{C}(\pi \Sigma), \mathbb{S}(s \vec{r} \vec{Z} \uplus \{z\}) \rightarrow}{\mathbb{C}(\pi \Sigma \uplus \sigma', \mathbb{S}(s \vec{r} \vec{Z}))}$		

Figure 5. The run-time system.

though we require that patterns have exact meets, in the sense that if  $z \sqsubseteq z_1$  and  $z \sqsubseteq z_2$  then  $z \sqsubseteq z_1 \sqcap z_2$ .

As with patterns, we leave the *actions*  $\gamma$  supported by the switches abstract, as different switches support different kinds of actions. However, actions must include  $\Omega$ , which sends packets to the controller. Moreover, although we do not require it formally, we expect that switches will support simple forward/modify actions. Many switches also support more sophisticated actions—e.g., Devoflow switches [18] can make conditional forwarding decisions based on whether their ports are up or down. Because switch actions  $\gamma$  are held abstract, we assume a function  $\text{Action}(\gamma, p)$  that executes an action  $\gamma$  on  $p$  and generates a set of transmissions  $T'$ . Typically  $T'$  will be a singleton set with a single transmission containing the packet being forwarded. However, if the policy dictates that the packet should be forwarded to multiple switches,  $T'$  will be bigger.

Given a packet  $p$  and a classifier  $\vec{r}$ , we match the packet against the classifier by finding the first rule whose pattern matches the packet. We write  $\vec{r} \rightsquigarrow^p z : \gamma$  for the matching judgment and  $\vec{r} \rightsquigarrow^p$  if  $p$  does not match any pattern in  $\vec{r}$ . More formally, we define classifier matching as follows; note that it selects the highest priority (leftmost) matching rule:

$$\frac{p \not\sqsubseteq z_1 \quad \dots \quad p \not\sqsubseteq z_{i-1} \quad p \sqsubseteq z_i}{(z_1 : \gamma_1, \dots, z_{i-1} : \gamma_{i-1}, z_i : \gamma_i, \dots, z_n : \gamma_n) \rightsquigarrow^p z_i : \gamma_i}$$

**The Molecular Machine.** Now we are ready to define the run-time system itself. We formalize its operational semantics as a *molecular machine*, in the style of Berry and Boudol’s chemical abstract machine [2]. The machine’s components, called *molecules*, are given on the left side of Figure 5.

The molecule  $\mathbb{C}(\pi|\Sigma)$  represents the controller machine running the NCore program  $\pi$  in state  $\Sigma$ . The molecule  $\mathbb{S}(s|\vec{r}|\vec{Z})$  represents switch  $s$  with packet classifier  $\vec{r}$  and local switch state  $\vec{Z}$ . The switch state records the patterns of rules that have been used to match packets but not yet queried and processed by the controller. Real switches use integer counters as state; for simplicity and elegance, we represent these integers in unary (using a multiset of patterns) in our formal system. A transmission molecule  $\mathbb{T}(s|p)$  represents a packet  $p$  en route to switch  $s$ . Finally, a help molecule  $\mathbb{H}(s|p)$  represents a request issued by switch  $s$  to the controller for assistance in processing packet  $p$ .

The machine’s operational semantics is defined on the right side of Figure 5. To lighten the notation, in this figure, we use the notation  $m_1, m_2$  as shorthand for  $\{m_1, m_2\}$ . Each operational rule may optionally be labelled with an *observation*  $o$ , which records when transmissions are processed.

The first three rules, E-SWITCHPROCESS, E-SWITCHHELP and E-SWITCHDROP, model the work done by switches to process

packets. There are three possibilities: either (1) the packet matches a rule with a non-controller action and the switch executes the action, (2) the packet matches a rule with a controller action and the switch generates a help molecule, or (3) the packet does not match any rule and is dropped.

The rule E-CONTROLLER models the work done by the controller to process help molecules. It first interprets the packet using its NCore program, generating new transmissions  $T'$  and new state  $\Sigma'$ . It then uses the additional information contained in the packet to construct new *reactive rules*  $\vec{r}'$  for the switch that generated the help molecule. These new rules typically allow the switch to process future packets locally, in the fast path, without requesting further help from the controller. We keep the definitions of the specific algorithms used to specialize the program and generate reactive rules abstract for now; they will be described in the next section.

The rule E-COLLECT models the work done by the controller to collect statistics about a particular switch-level rule. First, the switch picks a pattern  $z$  in the switch state to send to the controller. Next, the controller performs one step of processing using the query  $\phi$ .<sup>3</sup> Because switches do not store actual packets in their state, only patterns, the controller must construct a packet  $p$  that matches  $z$  to evaluate the query. We use the judgement  $\vec{r} \rightsquigarrow^p z : \gamma$  to represent fabrication of an appropriate packet and then appeal to the semantics of queries in the ordinary way. Finally, the controller stores the query result in its state. This rule places a subtle, but crucial constraint on classifiers: because we do not know which packet  $p'$  originally matched  $z$  on the switch, for correctness, the results of query evaluation must be the same no matter what packet  $p$  the controller chooses. To put it another way, if the rules placed on the switch are too broad and match too many different kinds of packets, then the results produced by E-COLLECT will be wrong.

## 5. The NCore Compiler

The NCore compiler performs two distinct tasks:

- *Classifier Generation:* given an NCore program, it builds a collection of classifiers, one for each switch in the network.
- *Reactive Specialization:* given a packet not handled by the current classifier installed on a switch, it generates additional rules that allow the switch to handle future packets with similar header fields efficiently, without consulting the controller.

<sup>3</sup>In practice, the controller would collect *all* of the statistics for a given rule at once; this can be modeled in our formal system by evaluating multiple E-COLLECT steps in sequence.

$\mathcal{S}(s, e) = \bar{z}$	$\mathcal{S}(s, h : \bar{w}) = \mathcal{O}^{\text{compile}}(h, \bar{w}), \top$ $\mathcal{S}(s, \text{inspect } i) = \top$ $\mathcal{S}(s, \text{switch } s') = \top$ $\mathcal{S}(s, e_1 \cap e_2) = \mathcal{S}(s, e_1) \times \mathcal{S}(s, e_2)$ $\mathcal{S}(s, \neg e) = \mathcal{S}(s, e)$
$\mathcal{S}(s, \tau) = \bar{z}$	$\mathcal{S}(s, e \rightarrow A) = \mathcal{S}(s, e)$ $\mathcal{S}(s, \tau_1 \cap \tau_2) = \mathcal{S}(s, \tau_1) \times \mathcal{S}(s, \tau_2)$ $\mathcal{S}(s, \neg \tau) = \mathcal{S}(s, \tau)$
$\mathcal{S}(s, \phi) = \bar{z}$	$\mathcal{S}(s, K \leftarrow e) = \mathcal{S}(s, e)$ $\mathcal{S}(s, \phi_1 \cap \phi_2) = \mathcal{S}(s, \phi_1) \times \mathcal{S}(s, \phi_2)$ $\mathcal{S}(s, \neg \phi) = \mathcal{S}(s, \phi)$
$\mathcal{S}(s, \pi) = \bar{z}$	$\mathcal{S}(s, (\tau, \phi)) = \mathcal{S}(s, \tau) \times \mathcal{S}(s, \phi)$

**Figure 6.** Skeleton Generation.

This section describes the key algorithms that implement these tasks. Due to space constraints, we focus on predicate compilation and, to a lesser extent, policy compilation. We discuss formal properties of the compiler in Section 6. Many of the definitions in this section are parameterized on abstract functions called *oracles*. Parameterizing the compiler in this way makes it independent of the details of the rapidly evolving SDN specifications.

### 5.1 Classifier Generation

The NCore classifier generator works in two phases. The first phase, *skeleton generation*, builds a sequence of switch patterns that form the backbone of the final classifier. The second, *action identification*, attaches a switch action  $\gamma$  to each pattern in the skeleton, using *Consistent Determinate Action Analysis* (CDAA) to select the action for each rule.

**Skeleton Generation.** A *skeleton* is half of a classifier that comprises a sequence of patterns. Intuitively, a “good” skeleton for a given policy is one where every pattern has a consistent and determinate action in the policy. Given a skeleton with this property, we can easily build a classifier by simply attaching the appropriate action to each pattern. If, however, the skeleton contains patterns for which the policy does not have a consistent and determinate action, then we must associate those patterns with the controller action  $\Omega$ , which causes packets to be diverted to the controller. This intuition about what makes a good skeleton cuts to the heart of why many OpenFlow programmers write programs using exact-match rules: Analyzing a classifier to identify consistent, determinate actions is easily done by hand when the classifier only contains exact-match rules. But when the classifier contains wildcard rules, the analysis becomes much more difficult. Thus, like register allocation, this analysis is better suited for computer processing than human processing.

Figure 6 presents the skeleton generation algorithm, which is formalized as a switch-indexed function  $\mathcal{S}(s, \cdot)$  from NCore programs to skeletons for  $s$ . One of the invariants of the algorithm is that it always generates a *complete skeleton*—one that matches all packets. Within this complete skeleton, the algorithm attempts to produce a pattern that matches the same set of packets as the predicate being compiled (relative to the patterns appearing before it in the classifier), and another set of patterns whose domain matches the complement of the predicate. While the algorithm attempts to generate skeletons with these properties, it does not always succeed, for two fundamental reasons: (1) it cannot analyze the decisions made by black-box inspectors—their results are indetermi-

nate, and (2) some predicates cannot be expressed practically (*i.e.*, without exponential blowup) in terms of switch patterns. For instance, building a good skeleton for a query that collects statistics on a per-source-IP (or per-destination-MAC, per-TCP port, *etc.*) basis would require generating distinct patterns for each of the approximately 4.3 billion IP addresses, which is clearly impractical. Nevertheless, in many common cases it successfully generates a good skeleton.

The first definition in Figure 6 generates skeletons for primitive predicates  $h : \bar{w}$  by invoking an oracle  $\mathcal{O}^{\text{compile}}$ . To handle cases where the switch hardware is more limited than these primitives,<sup>4</sup> we allow the skeleton to *overapproximate* the semantics of the primitive. However, when the skeleton is an overapproximation, the pattern in question will be associated with the  $\Omega$  action in the final classifier (the controller will resolve the pattern precisely). Note that the skeletons generated for primitives are terminated by  $\top$ . This establishes both properties mentioned above—completeness and the existence of complementary patterns.

The rule for compiling an inspector simply returns the  $\top$  pattern—it is impossible to analyze the properties of the black-box function to generate an accurate skeleton, so the rule overapproximates. Later in this section, we will see how reactive specialization allows us to overcome the initial ineffectiveness of this pattern.

The rule for compiling switch predicates is identical to the rule for inspectors, but the rationale is entirely different. If the switch name  $s'$  equals  $s$ , then the predicate switch  $s'$  matches all packets on  $s$  so the rule returns the  $\top$  pattern. If the two switch names are not equal, then the predicate matches no packets. The skeleton generation algorithm still returns  $\top$ —the pattern for the predicate itself would be  $\perp$ , with  $\top$  as the complement, but  $\perp$  may always be omitted as it does not match any packets.

The rule for intersections uses a cross-product operator, defined as follows, where  $\bigoplus_{i=1}^m$  denotes the sequence with elements indexed by  $i$  ranging from 1 to  $m$ :

$$(z_1, \dots, z_m) \times (z'_1, \dots, z'_n) = \bigoplus_{i=1}^m \bigoplus_{j=1}^n z_i \cap z'_j$$

One can think of the cross product as the “intersection” of two skeletons. We illustrate this property with an example, if the two patterns being combined are wildcard strings, the cross-product of the skeletons (1?, ??) and (10, 01, ??) is calculated as follows:

$$\begin{aligned} & (1? \cap 10, 1? \cap 01, 1? \cap ??, ?? \cap 10, ?? \cap 01, ?? \cap ??) \\ &= (10, \perp, 1?, 10, 01, ??) \\ &\equiv (10, 1?, 01, ??). \end{aligned}$$

Consider the relationship between the pattern 1? in the first skeleton and the pattern ?? in the second skeleton. The set of wildcard strings that simultaneously matches 1? and ?? in the context of their respective skeletons is {11} (note that 10 is matched in the second skeleton before one can reach ??). Likewise, the meet of 1? and ?? in the context of the cross product matches the same set. This property ensures that intersection preserves the “good” structure in skeletons.

Finally, the skeleton for a negated predicate  $\neg e$  is equal to the skeleton for  $e$ . This follows from the invariant that a generated skeleton must always be complete, and the goal that a skeleton should represent the semantics of both the predicate itself and its complement.

The cases for policies, queries, and programs are mostly straightforward. The skeleton generation algorithm ignores actions and keys but generate skeletons for the embedded predicates and surrounding set-theoretic operators.

<sup>4</sup> This is not uncommon—although bitwise patterns are convenient for programmers, they are not supported in many cases on OpenFlow.

$$\boxed{\mathcal{A}(s, \pi, \bar{z}) = \bar{r}} \\
\mathcal{A}(s, \pi, \cdot) = \cdot \\
\mathcal{A}(s, \pi, (\bar{z}, z)) = \begin{cases} \mathcal{A}(s, \pi, \bar{z}), z : \mathcal{O}^{\text{switch}}(A) \\ \text{if } s \mid \bar{z} \mid z \vdash \pi \rightsquigarrow (A, K) \\ \text{and } K \in \text{KeySet}[H] \\ \text{and } H \subseteq \text{headers}(z) \\ \mathcal{A}(s, \pi, \bar{z}), z : \Omega \\ \text{otherwise} \end{cases}$$

Figure 7. Action identification.

**Action Identification.** In the second phase of classifier generation, the compiler attaches an action to each pattern in the skeleton. Most of the work needed to do this is encapsulated in the algorithm for *Consistent, Determinate Action Analysis* (CDAA), which we represent using the following judgement:

$$s \mid \bar{z} \mid z \vdash \pi \rightsquigarrow (A, K).$$

If the above holds, then the procedure has determined that for every packet  $p$  that matches  $z$  in the context of  $\bar{z}$ , the program  $\pi$  associates the sets of actions  $A$  and keys  $K$  with  $p$  on switch  $s$ .

The CDAA is partial—for example, it cannot analyze black-box inspectors. Due to space limitations, we defer the full definition of the analysis, which is mostly straightforward, to the Appendix, along with the proofs of its formal properties. The key properties we have proven are (1) soundness, and (2) that the CDAA is sufficiently complete to ensure that reactive specialization makes progress; this is the main lemma used in the next section to prove the quiescence theorem.

One other ingredient we need is a switch action oracle ( $\mathcal{O}^{\text{switch}}$ ) that translates high-level actions into switch-level actions. For most actions, the translation is direct and obvious. For instance, it is typically the identity function for simple forwarding actions.

The action identification procedure is defined in Figure 7. It considers each pattern in the skeleton, and either (1) uses CDAA to find a consistent determinate action to associate with the pattern, checks that the pattern is sufficiently fine-grained for the queries, and creates an action using the switch action oracle, or (2) concludes that either a consistent determinate action does not exist (or cannot be found yet) or the pattern is insufficiently fine-grained for the queries, and associates the controller action with the pattern. To ensure correct implementation of queries we assume a function that computes a set of headers  $\text{headers}(z)$  from a pattern  $z$ . The key property of this function is that if  $h \in \text{headers}(z)$ , we know all packets  $p \sqsubseteq z$  have the same value in header  $h$ .

## 5.2 Reactive Specialization

The classifier generation algorithm described in the preceding section builds classifiers that can be installed on switches. But as we saw, it has some serious limitations. Many policies are dynamic and use black-box inspectors that cannot be analyzed. Even for purely static policies, a set of rules that matches all packets on the switches can be large—much larger than the actual space of flows that will be inhabited in a given network. To deal with these situations, we define *reactive specialization*, a powerful generalization of the simple, reactive, microflow-based strategy implemented manually by OpenFlow programmers. We define reactive specialization as the composition of three basic operations: (1) *program refinement*, which expands the program relative to a new packet witnessed by the controller, (2) *regeneration*, which translates the expanded program to a classifier using the tools of the previous sec-

$$\boxed{\mathcal{R}(x, e) = e'} \\
\mathcal{R}((\Sigma, s, p), h : \bar{w}) = h : \bar{w} \cup h : \mathcal{O}^{\text{refine}}(h, \bar{w}, p) \\
\mathcal{R}(x, \text{switch } s) = \text{switch } s \\
\mathcal{R}(x, \text{inspect } i) = \begin{cases} \text{inspect } i \cup e' & \text{if } \text{invariant}(x, i) \wedge i(x) \\ \text{inspect } i - e' & \text{if } \text{invariant}(x, i) \wedge \neg i(x) \\ \text{inspect } i & \text{if } \neg \text{invariant}(x, i) \end{cases} \\
\text{where } i : \text{Snapshot}[H] \rightarrow \text{Bool} \\ \text{and } e' = \text{similar}(x, H) \\
\mathcal{R}(x, e_1 \cap e_2) = \mathcal{R}(x, e_1) \cap \mathcal{R}(x, e_2) \\
\mathcal{R}(x, \neg e) = \neg \mathcal{R}(x, e)$$

$$\boxed{\mathcal{R}(x, \tau) = \tau'} \\
\mathcal{R}(x, e \rightarrow A) = (\mathcal{R}(x, e)) \rightarrow A \\
\mathcal{R}(x, \tau_1 \cap \tau_2) = \mathcal{R}(x, \tau_1) \cap \mathcal{R}(x, \tau_2) \\
\mathcal{R}(x, \neg \tau) = \neg \mathcal{R}(x, \tau)$$

$$\boxed{\mathcal{R}(x, \phi) = \phi'} \\
\mathcal{R}(x, K \leftarrow e) = K \leftarrow (e' \cup (e' \cap \text{similar}(x, H))) \\
\text{where } K \in \text{KeySet}[H] \\ \text{and } e' = \mathcal{R}(x, e) \\
\mathcal{R}(x, \phi_1 \cap \phi_2) = \mathcal{R}(x, \phi_1) \cap \mathcal{R}(x, \phi_2) \\
\mathcal{R}(x, \neg \phi) = \neg \mathcal{R}(x, \phi)$$

$$\boxed{\mathcal{R}(x, \pi) = \pi'} \\
\mathcal{R}(x, (\tau, \phi)) = (\mathcal{R}(x, \tau), \mathcal{R}(x, \phi))$$

Figure 8. NCore refinement.

tion, and (3) *reactive pruning*, which extracts new rules out of the regenerated classifier.

**Program Refinement.** When the controller receives a new packet that a switch could not handle, it interprets the program with respect to the packet, switch, and its current state. The idea in program refinement is to use that information to augment the program with additional structural information so that we can build a specialized classifier that handles future similar packets on the switch. Figure 8 defines the refinement function. The key invariant for this program transformation is that the semantics of the old and new programs are identical. However, syntactically, the new program will typically have a different structure, as the transformation uses the packet to unfold inspectors and queries. This makes skeleton generation more precise and the recompiled program more effective.

The rules for refining a predicate appear at the top of Figure 8. The first rule uses the refinement oracle  $\mathcal{O}^{\text{refine}}$  to refine basic predicates. Unlike the compilation oracle, which may *overapproximate* the predicate, the refinement oracle *underapproximates* it, so that the rest of the compilation infrastructure will be able to generate an implementable switch-level rule that matches the given packet. Importantly, because the new predicate is the union of the old predicate and an underapproximation, the overall semantics is unchanged. In some cases, especially if the switch-supported patterns are weak, the best underapproximation the refinement oracle can generate is an exact-match pattern. In many other cases, however, if the switch has prefix matching or wildcards, the refinement oracle will produce a rule that matches many more packets.

The second rule refines the switch predicate. Because the switch predicate already reveals the maximum amount of information, it cannot be refined further.

The rule for inspectors is the most interesting. Recall that if an inspector  $i$  is invariant with respect to a snapshot  $x$  (which includes the switch, packet, and controller state) then it will return the same boolean decision on similar packets in every future controller state. A packet is *similar* to another when it agrees on the headers  $H$  that the given inspector  $I$  examines. In such a case it is safe to augment the inspector with a predicate that is independent of the controller state, but matches similar packets in the same way as the inspector. The function  $\text{similar}(x, H)$  generates a predicate that matches packets similar to the packet in  $x$  on headers  $H$ :

$$\begin{aligned} \text{similar}((\Sigma, s, p), H) &= \text{switch } s \cap \bigcap_i e_i \\ &\text{where } e_i = h_i : p(h_i) \text{ for all } h_i \in H \end{aligned}$$

In the first clause for inspectors, if the inspector depends upon the headers  $H$  (as indicated by its indexed type), is invariant in the current state, and evaluates to true on the current snapshot  $(\Sigma, s, p)$ , then it may be safely refined by taking the union of the inspector and the similarity predicate  $\text{similar}((\Sigma, s, p), H)$ . The second inspector clause is similar except the inspector does *not* evaluate to true and hence the inspector predicate may be refined by subtracting out the similarity predicate. Finally, the third rule is used when the inspector is not invariant. In this case, no sound refinement exists—the boolean value returned by the inspector may change as the controller state changes—and so packets must still be diverted to the controller until the inspector stabilizes.

The other interesting element of refinement appears in the query translation. In particular, in the first rule, for queries  $K \leftarrow e$ , the algorithm expands the predicate  $e$  into  $e'$ , identifies the header set  $H$  relevant to  $K$  and generates the similarity predicate  $e''$ . The query predicate refinement is the union of  $e'$  (whose semantics, by induction, is identical to  $e$ ) and the intersection of  $e'$  and  $e''$ . The intersection identifies those packets similar to  $x$  that also satisfy the query predicate, and ensures that refinement preserves semantics.

**Reactive pruning.** In general, after the program has been refined and recompiled, some of the new rules will be useless—they will not process additional packets on the switch. We prune away these useless rules using a simple function  $\text{prune}(\bar{r}, p)$  that removes rules from  $\bar{r}$  that (1) send packets to the controller (adding such rules does not improve the efficiency of the switch), (2) have nothing to do with the packet  $p$  (meaning they are redundant rules not generated via specialization with respect to  $p$ ), or (3) overlap with a rule we removed earlier (to preserve the semantics of the rules).

### 5.3 Tying it all together

We now have the technical definitions needed to complete the definition of the NCore compiler and run-time system. In particular, we can now define (1)  $\text{GenClassifier}(\pi, s)$ , the classifier generator used to build the initial switch classifiers, and (2)  $\text{Specialize}(x, \pi)$ , the reactive specialist used by rule E-CONTROLLER to incrementally install new switch rules.

$$\frac{\begin{array}{l} \bar{z} = \mathcal{S}(s, \pi) \\ \bar{r} = \mathcal{A}(s, \pi, \bar{z}) \end{array}}{\text{GenClassifier}(\pi, s) = \bar{r}} \quad \frac{\begin{array}{l} \pi' = \mathcal{R}(x, \pi) \quad \bar{z} = \mathcal{S}(s, \pi') \\ \bar{r}' = \mathcal{A}(s, \pi', \bar{z}) \quad \bar{r}' = \text{prune}(\bar{r}, p) \end{array}}{\text{Specialize}(x, \pi) = \bar{r}'}$$

Finally, to initialize the molecular machine for NCore program  $\pi$  and switches  $s_1$  to  $s_n$ , we use the following functions:

$$\begin{aligned} \text{Init}(s, \pi) &= \mathbb{S}(s \mid \text{GenClassifier}(s, \pi) \mid \emptyset) \\ \text{Init}(\{s_1, \dots, s_n\}, \pi) &= \mathbb{C}(\pi \mid \emptyset, \text{Init}(s_1, \pi), \dots, \text{Init}(s_n, \pi)) \end{aligned}$$

## 6. Formal Properties

This section presents the essential elements of our two central theoretical results: (1) a proof of *functional correctness* for NCore, and (2) a proof of *quiescence*, another fundamental theorem which establishes that, when inspectors are invariant, the network eventually reaches a state in which all processing occurs efficiently on its switches.

**Functional Correctness** To demonstrate the correctness of the NCore compiler, we show that it inhabits the space between the asynchronous and synchronous reference machines. More formally, we prove that the asynchronous reference machine simulates the NCore molecular machine and the molecular machine simulates the synchronous reference machine. We state these properties formally below; the specific relations used to prove the simulations are given in the Appendix.

**Theorem 1** (Asynchronous weak simulation). The asynchronous machine  $(\pi, \emptyset, \bar{T}, \emptyset)$  weakly simulates the molecular machine  $(\text{Init}(S, \pi), \bar{T})$ .

**Theorem 2** (Synchronous weak simulation). The molecular machine  $(\text{Init}(S, \pi), \bar{T})$  weakly simulates the synchronous machine  $(\pi, \emptyset, \bar{T})$ .

As usual with simulation arguments, the key insight needed to complete the proof is finding the appropriate relations between the machines. The soundness of CDAA also plays a critical role.

**Quiescence** The quiescence theorem demonstrates that the NCore compiler effectively moves work off of the controller and onto switches, even when the program is expressed in terms of black-box inspectors. Formally, quiescence states that if all of the inspectors in the program are invariant, then the NCore compiler will eventually install rules on switches that handle all future traffic—*i.e.*, eventually, the system can reach a configuration where no additional packets need to be sent to the controller.

Before we can state the quiescence theorem precisely, we need a few definitions and supporting lemmas. Given a multiset of run-time system molecules  $\bar{M}$ , the *help multiset* of  $\bar{M}$ , written  $\text{Help}(\bar{M})$ , is the multiset of help molecules in  $\bar{M}$ . The *controller set* of  $\bar{M}$ , written  $\text{controllerPackets}(\bar{M})$  is the set of switch-packet pairs  $(s, p)$  that generate a help molecule when  $p$  is processed using the classifier currently installed on  $s$ . The first lemma, *controller set monotonicity*, states that the set of packets that require processing on the controller never increases:

**Lemma 1** (Controller set monotonicity). If  $(\text{Init}(S, \pi), \bar{T}) \rightarrow^* \bar{M}' \xrightarrow{o} \bar{M}''$  then  $\text{controllerPackets}(\bar{M}'') \subseteq \text{controllerPackets}(\bar{M}')$ .

As usual, the notation  $\bar{M} \rightarrow^* \bar{M}'$  denotes the reflexive, transitive closure of the single step judgement, ignoring observations.

Next we introduce two properties a program  $\pi$  must have to ensure quiescence: it must be *realizable* and *fully invariant*. Intuitively,  $\pi$  is realizable if it does not use features that cannot be implemented on switches—*e.g.*, in OpenFlow, matching on the payloads of packets or modifying multiple copies of a packet in different ways. Formally, we split realizability into two simpler properties: fully matchable and fully actionable. A program  $\pi$  is *fully matchable* if for every basic predicate  $h : \bar{w}$  appearing in it and every packet  $p$ , if  $\mathcal{O}^{\text{refine}}(h, w, p) = w'$  and  $\mathcal{O}^{\text{compile}}(h, w') = z$ , then for every snapshot  $(\Sigma, s, p')$ , we have  $p' \sqsubseteq z$  implies  $(\Sigma, s, p') \in \llbracket h : w \rrbracket$ . A program  $\pi$  is *fully actionable*, if, for all switches  $s$ , skeletons  $(\bar{z}, z)$ , and programs  $\pi'$ , if  $\llbracket \pi \rrbracket = \llbracket \pi' \rrbracket$  and  $s \mid \bar{z} \mid z \vdash \pi' \rightsquigarrow (A, K)$ , then  $\mathcal{O}^{\text{switch}}(A) \neq \Omega$ . Finally, a program  $\pi$  is *fully invariant* on  $\Sigma$  if, for every switch  $s$ , packet  $p$  and predicate  $\text{inspect } i$  appearing in it we have invariant  $((\Sigma, s, p), i)$ .

Now we are ready to prove the key lemma needed for quiescence. The *controller set progress* lemma states that if the controller program is realizable and has become fully invariant, then every time the controller processes a help molecule, the controller set becomes strictly smaller. In other words, every help molecule contains enough information (and the compiler is powerful enough to exploit it) for the E-CONTROLLER rule to generate useful new classifier rules. The full proof depends deeply on several completeness properties of CDAA; the Appendix contains further details.

**Lemma 2** (Controller set progress). For every realizable program  $\pi$ , if  $(\text{Init}(S, \pi), \bar{T}) \rightarrow^* (\mathbb{C}(\pi \mid \Sigma), \bar{M}') \xrightarrow{s,p} \bar{M}''$ , such that  $\pi$  is fully invariant on  $\Sigma$  and  $\text{Help}(\bar{M}'') \subset \text{Help}(\bar{M}')$ , then  $\text{controllerPackets}(\bar{M}'') \subset \text{controllerPackets}(\bar{M}')$ .

Quiescence follows as a simple corollary from controller set monotonicity and progress, as the total number of possible packets is finite. The precise statement of quiescence says that the run-time system *may* (as opposed to *does*) quiesce, because the machine may non-deterministically choose to continue forwarding packets using the switches instead of processing the remaining help molecules. Formally, a machine configuration  $\bar{M}$  *may quiesce* if there exists a configuration  $\bar{M}'$  such that  $\bar{M} \rightarrow^* \bar{M}'$  and the rule E-CONTROLLER cannot be used in any derivation in the operational semantics starting from  $\bar{M}'$ . With this definition in hand, we can state quiescence as follows:

**Corollary 1** (Quiescence). For every realizable program  $\pi$ , if  $(\text{Init}(S, \pi), \bar{T}) \rightarrow^* (\mathbb{C}(\pi \mid \Sigma), \bar{M}')$  and  $\pi$  is fully invariant on  $\Sigma$ , then  $\bar{M}'$  can quiesce.

## 7. Implementation and Evaluation

We have built a full working implementation of the NCore compiler in Haskell. It closely follows the formal definitions presented in this paper, modulo a few optimizations, which are described below. The core algorithms including skeleton generation, reactive specialization, and action identification are formulated in terms of type classes (e.g., lattices for switch patterns) and oracle functions. This design makes it easy to instantiate the compiler to target different switch architectures—one simply has to define instances for a few type classes and provide the appropriate oracles. We have built two back-ends, both based on OpenFlow switches. The first generates the coarse-grained rules from reactive specialization described in this paper. The other, used for comparison, generates exact-match rules, which emulates the microflow-based techniques commonly used by programmers today. The total amount of code needed to define these instances was only about 200 lines of code, compared to approximately 1600 lines for the compiler itself.

**Optimizations.** The implementation uses a number of heuristic optimizations to avoid the combinatorial blowup that would result from doing skeleton generation and reactive specialization naively. For example, it applies algebraic rewritings on-the-fly to remove useless patterns and rules and reduce the size of the intermediate patterns and classifiers it needs to manipulate. The skeleton generation algorithm identifies and removes patterns completely “shadowed” by other patterns and patterns whose effect is “covered” by a larger pattern lower down with the same actions. Finally, it uses a cheap static analysis to characterize the packets that will not be sent to the controller, which avoids invoking the more costly CDAA in many cases. Although these heuristics are simple, in our experience they go a long way toward ensuring reasonable performance.

**Evaluation.** To evaluate our implementation, we developed an instrumented version of the run-time system that collects statistics about the sizes of the classifiers and the amount of traffic handled

on switches (as opposed to being sent to the controller). Because space for classifiers is a limited resource on switches, and because the cost of diverting a packet to the controller slows down its processing by orders of magnitude, these metrics quantify some of the most critical aspects of the system.

We compared the performance of the “full” (which makes use of all OpenFlow rules, including wildcards) and “ $\mu$ flow” (which only generates exact-match rules) compilers on the following programs:

- **Static Policy (SP):** implements the simple static policy described at the end of Section 2. This benchmark measures the (in)efficiency of microflow-based compilation strategies.
- **Static Policy with Query (SPQ):** forwards packets using the same policy as in the SP test but also collects traffic statistics for each host. Due to the query, this program cannot be directly compiled to a switch classifier—at least, not without expanding all 4.3 billion possible hosts! Thus, this benchmark measures the efficiency of reactive compilation.
- **Black-Box Policy with Query (BPQ):** forwards packets and collects traffic statistics using the authentication application presented in Section 2. This benchmark measures the performance of a more realistic application implemented using both queries and black-box functions.

To drive these experiments, we generated packets using *fs* [23], a tool that synthesizes realistic packet traces from several statistical parameters. We ran each experiment on 100K packets in total. The results are shown in Figure ???. The graphs on the top row show the number of packets that “missed” and had to be sent to the controller against the total number of packets processed. Likewise, the graphs on the bottom row show the size of the compiled classifier, in terms of number of rules, versus total packets. The table at the right gives the final results after all 100K packets were processed.

In terms of the proportion of packets processed on switches, the full OpenFlow compiler outperforms the microflow-based compiler on all of the benchmarks. On the SP benchmark, the full compiler generates a classifier that completely handles the policy, so no packets are sent to the controller. (The line for the full compiler overlaps with the x-axis.) The microflow compiler, of course, diverts a packet to the controller for each distinct microflow. On the SPQ benchmark, after seeing a packet from each unique host involved in the query, the full compiler generates wildcard rules (via reactive specialization) that handle all future traffic from the host—many more packets than the exact-match rule produced by the microflow compiler. On this benchmark, it is worth noting that the classifiers produced by the full compiler are larger than the ones produced by the microflow compiler, especially initially. This is due to the fact that the full compiler sometimes generates multiple rules in response to a single controller packet, attempting to cover a broad space of future similar packets, whereas the microflow compiler predictably generates a single microflow for each controller packet. One can see that the work done by the full compiler pays off in terms of the number of packets that must be diverted to the controller. Moreover, over time, the size of the microflow compiler-generated classifier approaches that of the full compiler. Lastly, the BPQ experiment demonstrates that the full compiler generates more effective classifiers than the microflow compiler, even in the presence of black-box functions that it cannot analyze directly. Note that a large number of packets must be diverted to the controller in any correct implementation—at the start, the black-box is not invariant for any host. The difference between the two compilers starts to become clear toward the end of the experiment.

## 8. Related Work

Building on ideas first proposed in Ethane [4] and 4D [10], NOX [11] was the first concrete system to popularize what is currently known as software-defined networking. It provides an event-driven interface to OpenFlow [17] and requires programmers construct reactive programs manually out of callbacks and explicit, switch-level packet-processing rules. There are numerous examples of network applications built on top of NOX using microflows [12, 13, 28], but relatively few that use wildcard rules (though Wang’s load balancer [27] is a nice example of the latter).

Networking researchers are now actively developing next-generation controller platforms. Some of them, such as Beacon [1] (designed for Java) and Nettle [26] (designed for Haskell) provide elegant OpenFlow interfaces for new programming languages. Others, such as Onix [15], and Maestro [3] improve scalability and fault tolerance through parallelization and distribution. None of these systems automatically generate reactive protocols or provide formal semantics or correctness guarantees like NCore does.

NCore was inspired in part by Frenetic [9], a recently proposed high-level language for OpenFlow networks. NCore borrows from Frenetic the idea that network programs should be split into two parts: one for specification of forwarding policies and one for specification of queries. The major linguistic advance of NCore over Frenetic involves the inclusion of arbitrary functions in its policies and queries. In addition, NCore has a well-defined semantics, a parameterized, lattice-theoretic compilation strategy that uses the full capabilities of the switches, and proofs of correctness and quiescence of its implementation. Frenetic’s implementation uses inefficient microflows, is ad hoc, tailored to OpenFlow, and has no theoretical analysis of formal properties.

Both NCore and NDLog [16] use high-level languages to program networking infrastructure, but the similarities end there. NDLog programs are written in an explicitly distributed style whereas high-level NCore programs written as if the program has an omniscient, centralized view of the entire network. The NCore implementation automatically partitions work onto a distributed set of switches and synthesizes a reactive communication protocol that simulates the semantics of the high-level language.

Part of the job of the NCore compiler is to generate efficient packet classifiers. Most previous research in this area (see Taylor [25] for a survey) focuses on static compilation. The NCore compiler generates classifiers in the face of non-static policies, with unknown black box functions, and synthesizes a distributed switch-controller implementation. Bro [22], Snortan [7], Shangri-La [5] and FPL-3E [6] compile rich packet-filtering and monitoring programs, designed to secure networks and detecting intrusions, down to special packet-processing hardware and FPGAs. The main difference between NCore and all of these systems is that they are limited to a single device. They do not address the issue of how to program complex, dynamic policies for a collection of interconnected switches and they do not synthesize the distributed communication patterns between the switches and controller.

Active Networking, as in the SwitchWare project [24], shares many high-level goals with Software-Defined Networking, but the implementation strategy is entirely different. The former uses smart switches to interpret programs encapsulated in packets, while the latter uses dumb switches controlled by a remote, centralized host.

## References

- [1] Beacon: A java-based OpenFlow control platform., Nov 2010. See <http://www.beaconcontroller.net>.
- [2] G. Berry and G. Boudol. The chemical abstract machine. In *POPL*, pages 81–94, 1990.
- [3] Z. Cai, A. Cox, and T. Ng. Maestro: A system for scalable OpenFlow control. Technical Report TR10-08, Rice University, Dec 2010.
- [4] M. Casadod, M. Freedman, J. Pettit, J. Luo, N. Gude, N. McKeown, and S. Shenker. Rethinking enterprise network control. *Trans. on Networking.*, 17(4), Aug 2009.
- [5] M. Chen, X. Li, R. Lian, J. Lin, L. Liu, T. Liu, and R. Ju. Shangri-la: Achieving high performance from compiled network applications while enabling ease of programming. In *PLDI*, Jun 2005.
- [6] M. Cristea, C. Zissulescu, E. Deprettere, and H. Bos. FPL-3E: Towards language support for reconfigurable packet processing. In *SAMOS*, pages 201–212, Jul 2005.
- [7] S. Egorov and G. Savchuk. *SNORTRAN: An Optimizing Compiler for Snort Rules*. Fidelis Security Systems, 2002.
- [8] D. Erickson et al. A demonstration of virtual machine mobility in an OpenFlow network, Aug 2008. Demo at *ACM SIGCOMM*.
- [9] N. Foster, R. Harrison, M. Freedman, C. Monsanto, J. Rexford, A. Story, and D. Walker. Frenetic: A network programming language. In *ICFP*, Sep 2011.
- [10] A. Greenberg, G. Hjalmtysson, D. Maltz, A. Myers, J. Rexford, G. Xie, H. Yan, J. Zhan, and H. Zhang. A clean slate 4D approach to network control and management. *SIGCOMM Comput. Commun. Rev.*, 35:41–54, October 2005.
- [11] N. Gude, T. Koponen, J. Pettit, B. Pfaff, M. Casado, N. McKeown, and S. Shenker. NOX: Towards an operating system for networks. *SIGCOMM CCR*, 38(3), 2008.
- [12] N. Handigol, S. Seetharaman, M. Flajslik, N. McKeown, and R. Johari. Plug-n-Serve: Load-balancing web traffic using OpenFlow, Aug 2009. Demo at *ACM SIGCOMM*.
- [13] B. Heller, S. Seetharaman, P. Mahadevan, Y. Yiakoumis, P. Sharma, S. Banerjee, and N. McKeown. ElasticTree: Saving energy in data center networks. In *NSDI*, Apr 2010.
- [14] L. Jose, M. Yu, and J. Rexford. Online measurement of large traffic aggregates on commodity switches. In *Hot-ICE*, Mar 2011.
- [15] T. Koponen, M. Casado, N. Gude, J. Stribling, L. Poutievski, M. Zhu, R. Ramanathan, Y. Iwata, H. Inoue, T. Hama, and S. Shenker. Onix: A distributed control platform for large-scale production networks. In *OSDI*, Oct 2010.
- [16] B. Loo, J. Hellerstein, I. Stoica, and R. Ramakrishnan. Declarative routing: Extensible routing with declarative queries. In *SIGCOMM*, pages 289–300, 2005.
- [17] N. McKeown, T. Anderson, H. Balakrishnan, G. Parulkar, L. Peterson, J. Rexford, S. Shenker, and J. Turner. Openflow: Enabling innovation in campus networks. *SIGCOMM CCR*, 38(2):69–74, 2008.
- [18] J. Mogul, J. Tourrilhes, P. Yalagandula, P. Sharma, A. Curtis, and S. Banerjee. DevoFlow: Cost-effective flow management for high performance enterprise networks. In *HotNets*, pages 1:1–1:6, 2010.
- [19] A. Nayak, A. Reimers, N. Feamster, and R. Clark. Resonance: Dynamic access control in enterprise networks. In *WREN*, Aug 2009.
- [20] The Open Networking Foundation, Mar 2011. See <http://www.opennetworkingfoundation.org/>.
- [21] OpenFlow, Nov 2010. See <http://www.openflowswitch.org>.
- [22] Vern Paxson. Bro: A system for detecting network intruders in real-time. *Computer Networks*, 31(23–24):2435–2463, Dec 1999.
- [23] J. Sommers, R. Bowden, B. Eriksson, P. Barford, M. Roughan, and N. Duffield. Efficient network-wide flow record generation. In *INFOCOM*, pages 2363–2371, 2011.
- [24] SwitchWare. <http://www.cis.upenn.edu/~switchware>, 1997.
- [25] D. Taylor. Survey and taxonomy of packet classification techniques. *ACM Comput. Surv.*, 37:238–275, September 2005.
- [26] A. Voellmy and P. Hudak. Nettle: Functional reactive programming of OpenFlow networks. In *PADL*, Jan 2011.
- [27] R. Wang, D. Butnariu, and J. Rexford. OpenFlow-based server load balancing gone wild. In *Hot-ICE*, Mar 2011.
- [28] K. Yap, M. Kobayashi, R. Sherwood, T. Huang, M. Chan, N. Handigol, and N. McKeown. OpenRoads: Empowering research in mobile networks. *SIGCOMM Comput. Commun. Rev.*, 40(1):125–126, 2010.

## A. NCore Auxiliary Definitions and Properties

### A.1 Elided Skeleton and Classifier Properties

We define packet lookup  $\bar{z} \rightsquigarrow^p z$  in skeletons (and its complement failed lookup:  $\bar{z} \not\rightsquigarrow^p$ ) analogously to lookup in classifiers:

$$\frac{p \not\subseteq z_1 \quad \cdots \quad p \not\subseteq z_{i-1} \quad p \subseteq z_i}{(z_1, \dots, z_{i-1}, z_i, \dots, z_n) \rightsquigarrow^p z_i}$$

We define  $\text{packets}(\bar{z})$  to be the set of all packets  $p$  such that there exists a pattern  $z$  where  $\bar{z} \rightsquigarrow^p z$ . We define the  $\text{packets}(\bar{z}, z)$ , as  $\text{packets}(z) \setminus \text{packets}(\bar{z})$ , which captures the packets defined by  $z$  with respect to its prefix. Because classifiers are processed left to right,  $\text{packets}(\bar{z}, z)$  tells us the set of packets that will actually match  $z$  when preceded by  $\bar{z}$ .

**Lemma 3** (Packets matched by cross product). If  $\text{packets}(\bar{z}_1, z_1) = P$  and  $\text{packets}(\bar{z}_2, z_2) = P'$ , then for all  $\bar{z}$ ,

$$\text{packets}((\bar{z}_1 \times (\bar{z}_2, z_2, \bar{z}), z_1 \times \bar{z}_2), z_1 \cap z_2) = P \cap P'.$$

Given a classifier  $\bar{r}$ , we define  $\text{packets}(\bar{r})$  to be the set of packets  $p$  such that  $\bar{r} \rightsquigarrow^p z : \gamma$  for some  $z$  and  $\gamma$ . Define  $\text{controllerPackets}(\bar{r})$  to be the set of all packets  $p$  such that  $\bar{r} \rightsquigarrow^p z : \Omega$ .

### A.2 Oracles

This subsection states required properties of the oracles used directly in compilation. We use the type Wildcard for Vector WildcardBit.

$$\begin{aligned} \text{Compilation oracle } \mathcal{O}^{\text{compile}} &\in \text{Header} \times \text{Wildcard} \rightarrow \text{Pattern} \\ \text{Expansion oracle } \mathcal{O}^{\text{refine}} &\in \text{Header} \times \text{Wildcard} \times \text{Packet} \rightarrow \text{Wildcard} \\ \text{Switch oracle } \mathcal{O}^{\text{switch}} &\in \text{Action} \rightarrow \text{SwitchAct} \end{aligned}$$

**Definition 1** (Compilation oracle correctness). If  $\mathcal{O}^{\text{compile}}(h, w) = z$ , then for all  $(\Sigma, s, p) \in \llbracket h : w \rrbracket$ , we have  $p \subseteq z$ .

**Definition 2** (Expansion oracle correctness). If  $\mathcal{O}^{\text{refine}}(h, w, p) = w'$ , then

- $\llbracket h : w' \rrbracket \subseteq \llbracket h : w \rrbracket$  and
- if there exists a  $w''$  such that  $\mathcal{O}^{\text{compile}}(h, w'') = z$  and  $p \subseteq z$  only if  $(\Sigma, s, p) \in \llbracket h : w \rrbracket$ , then  $w' = w''$ .

**Definition 3** (Switch oracle correctness). If  $\mathcal{O}^{\text{switch}}(A) = \gamma$  and  $\gamma \neq \Omega$ , then  $\text{Action}(\gamma, p) = T'$  if and only if  $\llbracket A \rrbracket(p) = T'$ .

### A.3 Pruning

There are a variety of different prunes that can be done. Regardless of the prune chosen, it must satisfy the following correctness conditions.

**Definition 4** (Prune correctness). Let  $\text{prune}(\bar{r}, p) = \bar{r}'$ . Then we have that:

- If  $\bar{r}$  on switch  $s$  is correct with respect to  $\pi$  then  $(\bar{r}', \bar{r})$  on  $s$  is correct with respect to  $\pi$ ,
- $\text{controllerPackets}(\bar{r}') = \emptyset$ ,
- If  $p \in \text{packets}(\bar{r}) \setminus \text{controllerPackets}(\bar{r})$ , then  $p \in \text{packets}(\bar{r})$ .

### A.4 Consistent Definitive Action Analysis

Consistent definitive action analysis (CDAA) is an analysis that attempts to determine the consistent set of actions and keys associated with a pattern. To be more specific, it determines whether a predicate  $e'$  is *absorbed* or *annihilated* by a pattern  $z$ . Notice that there are no rules involving inspectors – the analysis can not determine the actions produced by an inspector. If it is *necessary* for the analysis to understand the semantics of an inspector, it will fail to deduce a consistent action.

The analysis uses one further oracle, the CDAA oracle.

$$\text{CDAA oracle } \mathcal{O}^{\text{CDAA}} \in \text{Pattern} \times \text{Header} \rightarrow \text{Wildcard}$$

**Definition 5** (Disjoint bitvectors). We say  $\bar{w}$  is disjoint from  $\bar{w}'$  when there does not exist a bitvector  $\bar{b}$  such that  $\bar{b} \subseteq \bar{w}$  and  $\bar{b} \subseteq \bar{w}'$

**Definition 6** (Static analysis oracle correctness).  $p \subseteq z$  if and only if for every  $h$ , we have  $p(h) \subseteq \mathcal{O}^{\text{CDAA}}(z, h)$ .

$s \mid \bar{z} \mid z \vdash e$  absorbed,  $s \mid \bar{z} \mid z \vdash e$  annihilated

$$\frac{\text{AB-PRSMALLER} \quad \mathcal{O}^{\text{CDAA}}(z, h) = \bar{w}' \quad \bar{w}' \sqsubseteq \bar{w}}{s \mid \bar{z} \mid z \vdash h : \bar{w} \text{ absorbed}}$$

$$\frac{\text{AN-PRDISJOINT} \quad \mathcal{O}^{\text{CDAA}}(z, h) = \bar{w}' \quad \bar{w} \text{ and } \bar{w}' \text{ disjoint}}{s \mid \bar{z} \mid z \vdash h : \bar{w} \text{ annihilated}}$$

$$\frac{\text{AN-VOID} \quad \mathcal{O}^{\text{CDAA}}(z, h) = \bar{w}' \quad \mathcal{O}^{\text{CDAA}}(z', h) = \bar{w}'' \quad \bar{w} \sqcap \bar{w}' \sqsubseteq \bar{w} \sqcap \bar{w}''}{s \mid \bar{z}, z', \bar{z}' \mid z \vdash h : \bar{w} \text{ annihilated}}$$

$$\frac{\text{AB-PRSWITCH}}{s \mid \bar{z} \mid z \vdash \text{switch } s \text{ absorbed}}$$

$$\frac{\text{AN-PRSWITCH} \quad s \neq s'}{s \mid \bar{z} \mid z \vdash \text{switch } s' \text{ annihilated}}$$

$$\frac{\text{AB-PRINTERSECT} \quad s \mid \bar{z} \mid z \vdash e_1 \text{ absorbed} \quad s \mid \bar{z} \mid z \vdash e_2 \text{ absorbed}}{s \mid \bar{z} \mid z \vdash e_1 \cap e_2 \text{ absorbed}}$$

$$\frac{\text{AN-PRINTERSECT-1} \quad s \mid \bar{z} \mid z \vdash e_1 \text{ annihilated}}{s \mid \bar{z} \mid z \vdash e_1 \cap e_2 \text{ annihilated}}$$

$$\frac{\text{AN-PRINTERSECT-2} \quad s \mid \bar{z} \mid z \vdash e_2 \text{ annihilated}}{s \mid \bar{z} \mid z \vdash e_1 \cap e_2 \text{ annihilated}}$$

$$\frac{\text{AB-PRNEG} \quad s \mid \bar{z} \mid z \vdash e \text{ annihilated}}{s \mid \bar{z} \mid z \vdash \neg e \text{ absorbed}}$$

$$\frac{\text{AN-PRNEG} \quad s \mid \bar{z} \mid z \vdash e \text{ absorbed}}{s \mid \bar{z} \mid z \vdash \neg e \text{ annihilated}}$$

$s \mid \bar{z} \mid z \vdash \tau \rightsquigarrow A$

$$\frac{\text{CDAA-POABPR} \quad s \mid \bar{z} \mid z \vdash e \text{ absorbed}}{s \mid \bar{z} \mid z \vdash (e \rightarrow A) \rightsquigarrow A}$$

$$\frac{\text{CDAA-POANPR} \quad s \mid \bar{z} \mid z \vdash e \text{ annihilated}}{s \mid \bar{z} \mid z \vdash (e \rightarrow A) \rightsquigarrow \emptyset}$$

$$\frac{\text{CDAA-POINTERSECT} \quad s \mid \bar{z} \mid z \vdash \tau \rightsquigarrow A \quad s \mid \bar{z} \mid z \vdash \tau' \rightsquigarrow A'}{s \mid \bar{z} \mid z \vdash (\tau \cap \tau') \rightsquigarrow A \cap A'}$$

$$\frac{\text{CDAA-POINTERSECTEMPTY-1} \quad s \mid \bar{z} \mid z \vdash \tau \rightsquigarrow \emptyset}{s \mid \bar{z} \mid z \vdash (\tau \cap \tau') \rightsquigarrow \emptyset}$$

$$\frac{\text{CDAA-POINTERSECTEMPTY-2} \quad s \mid \bar{z} \mid z \vdash \tau' \rightsquigarrow \emptyset}{s \mid \bar{z} \mid z \vdash (\tau \cap \tau') \rightsquigarrow \emptyset}$$

$$\frac{\text{CDAA-PONEG} \quad s \mid \bar{z} \mid z \vdash \tau \rightsquigarrow A}{s \mid \bar{z} \mid z \vdash \neg \tau \rightsquigarrow \neg A}$$

$s \mid \bar{z} \mid z \vdash \phi \rightsquigarrow K$

$$\frac{\text{CDAA-QABPR} \quad s \mid \bar{z} \mid z \vdash e \text{ absorbed}}{s \mid \bar{z} \mid z \vdash (K \leftarrow e) \rightsquigarrow K}$$

$$\frac{\text{CDAA-QANPR} \quad s \mid \bar{z} \mid z \vdash e \text{ annihilated}}{s \mid \bar{z} \mid z \vdash (K \leftarrow e) \rightsquigarrow \emptyset}$$

$$\frac{\text{CDAA-QINTERSECT} \quad s \mid \bar{z} \mid z \vdash \phi \rightsquigarrow K \quad s \mid \bar{z} \mid z \vdash \phi' \rightsquigarrow K'}{s \mid \bar{z} \mid z \vdash (\phi \cap \phi') \rightsquigarrow K \cap K'}$$

$$\frac{\text{CDAA-QINTERSECTEMPTY-1} \quad s \mid \bar{z} \mid z \vdash \phi \rightsquigarrow \emptyset}{s \mid \bar{z} \mid z \vdash (\phi \cap \phi') \rightsquigarrow \emptyset}$$

$$\frac{\text{CDAA-QINTERSECTEMPTY-2} \quad s \mid \bar{z} \mid z \vdash \phi \rightsquigarrow \emptyset}{s \mid \bar{z} \mid z \vdash (\phi \cap \phi') \rightsquigarrow \emptyset}$$

$$\frac{\text{CDAA-QNEG} \quad s \mid \bar{z} \mid z \vdash \phi \rightsquigarrow K}{s \mid \bar{z} \mid z \vdash \neg \phi \rightsquigarrow \neg K}$$

$s \mid \bar{z} \mid z \vdash \pi \rightsquigarrow (A, K)$

$$\frac{\text{CDAA-PROGRAM} \quad s \mid \bar{z} \mid z \vdash \tau \rightsquigarrow A \quad s \mid \bar{z} \mid z \vdash \phi \rightsquigarrow K}{s \mid \bar{z} \mid z \vdash (\tau, \phi) \rightsquigarrow (A, K)}$$

**Lemma 4** (Soundness of CDAA).

- If  $s \mid \bar{z} \mid z \vdash e$  absorbed, then for all  $(\Sigma, s, p)$  such that  $p \in \text{packets}(\bar{z}, z)$ , we have that  $(\Sigma, s, p) \in \llbracket e \rrbracket$ .
- If  $s \mid \bar{z} \mid z \vdash e$  annihilated, then for all  $(\Sigma, s, p)$  such that  $p \in \text{packets}(\bar{z}, z)$ , we have that  $(\Sigma, s, p) \notin \llbracket e \rrbracket$ .

*Proof.* Induction on the derivation of the hypotheses. □

**Theorem 3** (Soundness). If  $s \mid \bar{z} \mid z \vdash \pi \rightsquigarrow (A, K)$ , then for all  $(\Sigma, s, p)$  such that  $p \in \text{packets}(\bar{z}, z)$ , we have  $\llbracket \pi \rrbracket(\Sigma, s, p) = (\llbracket A \rrbracket(\Sigma, s, p), \llbracket K \rrbracket(\Sigma, s, p))$ .

## A.5 Elided Correctness Properties

**Definition 7** (Classifier correct with respect to policy). We say that classifier  $\bar{r}$  on switch  $s$  is correct with respect to policy  $\tau$  if for every  $p$  such that  $\bar{r} \rightsquigarrow^p z : \gamma$  and  $\gamma \neq \Omega$ , then for all  $\Sigma$ , we have  $\text{Action}(\gamma, p) = T'$  if and only if  $\llbracket \tau \rrbracket(\Sigma, s, p) = T'$ .

**Definition 8** (Classifier correct with respect to query). A switch classifier  $\bar{r}$  on switch  $s$  is correct with respect to  $\phi$  if for all  $z$  there exists a  $\Sigma'$  such that for all  $p$  such that  $\bar{r} \rightsquigarrow^p z : \gamma$  and  $\gamma \neq \Omega$ , then  $\llbracket \phi \rrbracket(\Sigma, s, p) = \Sigma'$ .

**Definition 9** (Classifier correct with respect to program). A switch classifier  $\bar{r}$  on switch  $s$  is correct with respect to  $\pi = (\tau, \phi)$  if it is correct with respect to  $\tau$  and correct with respect to  $\phi$ .

**Theorem 4** (Correctness of classifier generation). If  $\text{GenClassifier}(\pi, s) = \bar{r}$ , then  $\bar{r}$  on switch  $s$  is correct with respect to  $\pi$ .

*Proof.* Follows from soundness of CDAA and the axioms of the switch oracle. □

**Theorem 5** (Expansion correctness). For all  $x$  and  $\pi$ , we have that  $\llbracket \mathcal{R}(x, \pi) \rrbracket = \llbracket \pi \rrbracket$ .

*Proof.* By induction on  $\pi$ . □

Using the above lemmas and theorems, one can prove the following relations, establishing our functional correctness theorem.

**Definition 10** (Asynchronous simulation relation). Let

$$\begin{aligned} \bar{M} &= \mathbb{C}(\pi \mid \Sigma), \mathbb{S}(s_1 \mid \bar{r}_1 \mid \bar{Z}_1), \dots, \mathbb{S}(s_n \mid \bar{r}_n \mid \bar{Z}_n), \bar{L}, \bar{T} \\ M_{\text{async}} &= (\pi', \Sigma', \bar{T}_1, \bar{T}_2). \end{aligned}$$

Then  $\bar{M} \sim_{\text{async}} M_{\text{async}}$  if

- $\llbracket \pi \rrbracket = \llbracket \pi' \rrbracket$
- For all  $i$ , classifier  $\bar{r}_i$  on switch  $s_i$  is correct with respect to  $\pi$ .
- $\Sigma = \Sigma'$
- There exists a bijection between  $\bar{T}_1$  and  $\bar{L} \uplus \bar{T}$ , where every mapping preserves  $(s, p)$ .
- There exists a bijection between transmissions  $\mathbb{T}(s \mid p) \in \bar{T}_2$  and patterns  $z \in \bar{Z}_i$ , where  $s_i = s$  and  $p \sqsubseteq z$ .

**Definition 11** (Synchronous simulation relation). Let

$$\begin{aligned} \bar{M} &= \mathbb{C}(\pi \mid \Sigma), \mathbb{S}(s_1 \mid \bar{r}_1 \mid \bar{Z}_1), \dots, \mathbb{S}(s_n \mid \bar{r}_n \mid \bar{Z}_n), \bar{L}, \bar{T} \\ M_{\text{sync}} &= (\pi', \Sigma', \bar{T}'). \end{aligned}$$

Then  $\bar{M} \sim_{\text{sync}} M_{\text{sync}}$  if

- $\llbracket \pi \rrbracket = \llbracket \pi' \rrbracket$
- For all  $i$ , classifier  $\bar{r}_i$  on switch  $s_i$  is correct with respect to  $\pi$ .
- $\Sigma = \Sigma'$
- There exists a bijection between  $\bar{L} \uplus \bar{T}$  and  $\bar{T}'$ , where every mapping preserves  $(s, p)$ .
- Each  $\bar{Z}_i$  is empty.

## A.6 Elided Quiescence Properties

**Definition 12** (Ideal packets matched on switch). Given a predicate  $e$ , define the set of packets ideally matched on switch  $s$  (written  $\text{packetsOn}(s, e)$ ) as all  $p$  such that  $(\Sigma, s, p) \in \llbracket e \rrbracket$  implies  $(\Sigma', s, p) \in \llbracket e \rrbracket$ , for all  $\Sigma$  and  $\Sigma'$ . Analogously, define the set of packets ideally dropped on the switch  $s$  (written  $\text{nPacketsOn}(s, e)$ ) as all  $p$  such that  $(\Sigma, s, p) \in \neg \llbracket e \rrbracket$  implies  $(\Sigma', s, p) \in \neg \llbracket e \rrbracket$ , for all  $\Sigma$  and  $\Sigma'$ .

**Definition 13** (anti-kernel). Define the anti-kernel with respect to  $(\Sigma, s)$  of a policy  $\tau$  (resp. a query  $\phi$ ), denoted  $\text{aker}(\Sigma, s, \tau)$  (resp.  $\text{aker}(\Sigma, s, \phi)$ ) as the set of all packets  $p$  such that  $\llbracket \tau \rrbracket(\Sigma, s, p) \neq \emptyset$  (resp.  $\llbracket \phi \rrbracket(\Sigma, s, p) \neq \emptyset$ ).

The following function predicts how many packets will go to the controller.

$\text{Inconsistent}(s, e) = P$

$$\frac{\text{INCON-PR EXACT} \quad \mathcal{O}^{\text{compile}}(h, \bar{w}) = z \quad \llbracket h : \bar{w} \rrbracket = \text{packets}(z)}{\text{Inconsistent}(s, h : \bar{w}) = \emptyset}$$

$$\frac{\text{INCON-PR OVERAPPROX} \quad \mathcal{O}^{\text{compile}}(h, \bar{w}) = z \quad \llbracket h : \bar{w} \rrbracket \neq \text{packets}(z)}{\text{Inconsistent}(s, h : \bar{w}) = \text{packets}(z)}$$

$$\frac{\text{INCON-PR WITCH}}{\text{Inconsistent}(s, \text{switch } s') = \emptyset}$$

$$\frac{\text{INCON-PR INSPECT}}{\text{Inconsistent}(s, \text{inspect } i) = \top}$$

$$\frac{\text{INCON-PR INTERSECT} \quad \text{Inconsistent}(s, e) = P \quad \text{Inconsistent}(s, e') = P'}{\text{Inconsistent}(s, e \cap e') = (P \cup P') \cap (\text{packetsOn}(s, e) \cup P) \cap (\text{packetsOn}(s, e') \cup P')}$$

$$\frac{\text{INCON-PR NEG} \quad \text{Inconsistent}(s, e) = P}{\text{Inconsistent}(s, \neg e) = P}$$

$$\boxed{\text{Inconsistent}(s, \tau) = P}$$

$$\frac{\text{INCON-POPR} \quad \text{Inconsistent}(s, e) = P}{\text{Inconsistent}(s, e \rightarrow A) = P}$$

$$\frac{\text{INCON-POINTERSECT} \quad \text{Inconsistent}(s, \tau) = P \quad \text{Inconsistent}(s, \tau') = P'}{\text{Inconsistent}(s, \tau \cap \tau') = (P \cup P') \cap (\text{aker}(\Sigma, s, \tau) \cup P) \cap (\text{aker}(\Sigma, s, \tau') \cup P')}$$

$$\frac{\text{INCON-PONEG} \quad \text{Inconsistent}(s, \tau) = P}{\text{Inconsistent}(s, \neg\tau) = P}$$

$$\boxed{\text{Inconsistent}(s, \phi) = P}$$

$$\frac{\text{INCON-QPR} \quad \text{Inconsistent}(s, e) = P}{\text{Inconsistent}(s, K \leftarrow e) = P}$$

$$\frac{\text{INCON-QINTERSECT} \quad \text{Inconsistent}(s, \phi) = P \quad \text{Inconsistent}(s, \phi') = P'}{\text{Inconsistent}(s, \phi \cap \phi') = (P \cup P') \cap (\text{aker}(\Sigma, s, \phi) \cup P) \cap (\text{aker}(\Sigma, s, \phi') \cup P')}$$

$$\frac{\text{INCON-QNEG} \quad \text{Inconsistent}(s, \phi) = P}{\text{Inconsistent}(s, \neg\phi) = P}$$

$$\boxed{\text{Inconsistent}(s, F) = P}$$

$$\frac{\text{INCON-PROGRAM} \quad \text{Inconsistent}(s, \tau) = P \quad \text{Inconsistent}(s, \phi) = P'}{\text{Inconsistent}(s, (\tau, \phi)) = (P \cup P') \cap (\text{aker}(\Sigma, s, \tau) \cup P) \cap (\text{aker}(\Sigma, s, \phi) \cup P')}$$

**Lemma 5** (Prefix extension).

- If  $s \mid \bar{z} \mid z \vdash e$  absorbed and  $\bar{z}'$  is a supersequence of  $\bar{z}$ , then  $s \mid \bar{z}' \mid z \vdash e$  absorbed.
- If  $s \mid \bar{z} \mid z \vdash e$  annihilated and  $\bar{z}'$  is a supersequence of  $\bar{z}$ , then  $s \mid \bar{z}' \mid z \vdash e$  annihilated.

*Proof.* Induction on the derivation of the hypotheses, ignoring the other members of  $\bar{z}'$  aside from  $\bar{z}$ . □

**Lemma 6** (Uniform meet). If  $\mathcal{S}(s, e) = (\bar{z}, z, \bar{z}')$ , then

- If  $s \mid \bar{z} \mid z \vdash e$  absorbed, then  $s \mid \bar{z} \times z' \mid z \cap z' \vdash e$  absorbed.
- If  $s \mid \bar{z} \mid z \vdash e$  annihilated, then  $s \mid \bar{z} \times z' \mid z \cap z' \vdash e$  annihilated.

*Proof.* Induction on the derivation of the hypotheses. □

**Lemma 7** (Cross product relations). For all  $\bar{z}$ , if  $\mathcal{S}(s, e_1) = (\bar{z}_1, z_1, \bar{z}'_1)$  and  $\mathcal{S}(s, e_2) = (\bar{z}_2, z_2, \bar{z}'_2)$  and  $\bar{z}' = \bar{z}_1 \times (\bar{z}_2, z_2, \bar{z})$ ,  $z_1 \times \bar{z}_2$ , then

- If  $s \mid \bar{z}_1 \mid z_1 \vdash e_1$  absorbed and  $s \mid \bar{z}_2 \mid z_2 \vdash e_2$  absorbed then  $s \mid \bar{z}' \mid z_1 \cap z_2 \vdash e_1 \cap e_2$  absorbed.
- If  $s \mid \bar{z}_1 \mid z_1 \vdash e_1$  annihilated or  $s \mid \bar{z}_2 \mid z_2 \vdash e_2$  annihilated then  $s \mid \bar{z}' \mid z_1 \cap z_2 \vdash e_1 \cap e_2$  annihilated.

*Proof.* Follows from uniform meet, prefix extension, and Ab-PrIntersect/An-PrIntersect. □

**Lemma 8** (Consistent predicate skeleton and absorption/annihilation). If  $\mathcal{S}(s, e) = (\bar{z}, z, \bar{z}')$  and  $\text{Inconsistent}(s, e) = P$ , then, for all  $\Sigma$ ,

- if  $\text{packets}(\bar{z}, z) \subseteq \text{packetsOn}(s, e) \setminus P$ , we have that  $s \mid \bar{z} \mid z \vdash e$  absorbed.
- if  $\text{packets}(\bar{z}, z) \subseteq \text{nPacketsOn}(s, e) \setminus P$ , we have that  $s \mid \bar{z} \mid z \vdash e$  annihilated.

*Proof.* Induction on  $e$ , using the cross-product relations and the algebraic properties of  $\text{Inconsistent}(s, e)$  to discharge cases in which we don't have useful induction hypotheses. □

**Theorem 6** (Consistent skeleton and CDAA). Let  $\mathcal{S}(s, \pi) = (\bar{z}, z, \bar{z}')$  and  $\text{Inconsistent}(s, \pi) = P$ . Then if  $\text{packets}(\bar{z}, z)$  disjoint from  $P$ , we have that  $s \mid \bar{z} \mid z \vdash \pi \rightsquigarrow (A, K)$ .

$$\text{HeadersAvail}(e, h) = P$$

$$\begin{aligned} \text{HeadersAvail}(h : w, h) &= \begin{cases} \text{packets}(z) & \text{if } \mathcal{O}^{\text{compile}}(h, w) = z \text{ and } \mathcal{O}^{\text{CDA}}(z, h) = \bar{b} \\ \emptyset & \text{otherwise} \end{cases} \\ \text{HeadersAvail}(\text{switch } s, h) &= \emptyset \\ \text{HeadersAvail}(\text{inspect } i, h) &= \emptyset \\ \text{HeadersAvail}(e_1 \cap e_2, h) &= \text{HeadersAvail}(e_1, h) \cup \text{HeadersAvail}(e_2, h) \\ \text{HeadersAvail}(\neg e, h) &= \text{HeadersAvail}(e, h) \end{aligned}$$

$$\text{HeadersAvail}(\tau, h) = P$$

$$\begin{aligned} \text{HeadersAvail}(e \rightarrow A, h) &= \text{HeadersAvail}(e, h) \\ \text{HeadersAvail}(\tau_1 \cap \tau_2, h) &= \text{HeadersAvail}(\tau_1, h) \cup \text{HeadersAvail}(\tau_2, h) \\ \text{HeadersAvail}(\neg \tau, h) &= \text{HeadersAvail}(\tau, h) \end{aligned}$$

$$\text{HeadersAvail}(\phi, h) = P$$

$$\begin{aligned} \text{HeadersAvail}(K \leftarrow e, h) &= \text{HeadersAvail}(e, h) \\ \text{HeadersAvail}(\phi_1 \cap \phi_2, h) &= \text{HeadersAvail}(\phi_1, h) \cup \text{HeadersAvail}(\phi_2, h) \\ \text{HeadersAvail}(\neg \phi, h) &= \text{HeadersAvail}(\phi, h) \end{aligned}$$

$$\text{HeadersAvail}(\pi, h) = P$$

$$\text{HeadersAvail}((\tau, \phi), h) = \text{HeadersAvail}(\tau, h) \cup \text{HeadersAvail}(\phi, h)$$

**Theorem 7** (Available headers of skeleton). Let  $\mathcal{S}(s, \pi) = (\bar{z}, z, \bar{z}')$ . Then there exists a  $\bar{b}$  for every  $p \in \text{packets}(\bar{z}, z)$  such that  $p(h) = \bar{b}$  if  $p \in \text{HeadersAvail}(\pi, h)$ .

**Definition 14** (Headers needed on snapshot). Let  $p \in \text{HeadersNeeded}(s, \pi, h)$  if for all  $\Sigma$ , we have  $\llbracket \pi \rrbracket(\Sigma, s, p) = (T', \sigma')$  and there exists  $(k, sp) \in \sigma'$  such that  $h \in \text{dom}(sp)$ .

We now define which packets we expect to see at the controller.

**Definition 15.** We let  $p \in \text{ToController}(s, \pi)$  if:

- $p \in \text{Inconsistent}(s, \pi)$  or
- there exists an  $h$  and a  $p$  such that  $p \in \text{HeadersNeeded}(s, \pi, h)$  and  $p \notin \text{HeadersAvail}(\pi, h)$ .

**Theorem 8** (Controller actions bounded). If  $\pi$  is realizable on  $s$ , then  $\mathcal{S}(s, \pi) = \bar{z}$  and  $\mathcal{A}(s, \pi, \bar{z}) = \bar{r}$  then  $\text{controllerPackets}(\bar{r}) \subseteq \text{ToController}(s, \pi)$ .

*Proof.* Induction on  $\pi$ . □

**Theorem 9** (Expansion monotonic). For all  $s, x$  and  $\pi$ , we have that  $\text{ToController}(s, \mathcal{R}(x, \pi)) \subseteq \text{ToController}(s, \pi)$ .

*Proof.* Induction on  $\pi$ . □

**Theorem 10** (Expansion invariance). For all  $(\Sigma, s, p)$ , if  $\pi$  is fully invariant and fully matchable, we have that  $p \notin \text{ToController}(s, \mathcal{R}((\Sigma, s, p), \pi))$ .

*Proof.* Induction on  $\pi$ . □