Machine-Verified Network Controllers

Arjun Guha
Cornell University
arjun@cs.cornell.edu

Mark Reitblatt
Cornell University
reitblatt@cs.cornell.edu

Nate Foster
Cornell University
jnfoster@cs.cornell.edu

Abstract

In many areas of computing, techniques ranging from testing to formal modeling to full-blown verification have been successfully used to help programmers build reliable systems. But although networks are critical infrastructure, they have largely resisted analysis using formal techniques. Software-defined networking (SDN) is a new network architecture that has the potential to provide a foundation for network reasoning, by standardizing the interfaces used to express network programs and giving them a precise semantics.

This paper describes the design and implementation of the first machine-verified SDN controller. Starting from the foundations, we develop a detailed operational model for OpenFlow (the most popular SDN platform) and formalize it in the Coq proof assistant. We then use this model to develop a verified compiler and run-time system for a high-level network programming language. We identify bugs in existing languages and tools built without formal foundations, and prove that these bugs are absent from our system. Finally, we describe our prototype implementation and our experiences using it to build practical applications.

Categories and Subject Descriptors F.3.1 [Specifying and Verifying and Reasoning about Programs]: Mechanical verification

Keywords Software-defined networking, OpenFlow, formal verification, Coq, domain-specific languages, NetCore, Frenetic.

1. Introduction

Networks are some of the most critical infrastructure in modern society and also some of the most fragile! Networks fail with alarming frequency, often due to simple misconfigurations or software bugs [8, 19, 30]. The recent news headlines contain numerous examples of network failures leading to disruptions: a configuration error during routine maintenance at Amazon triggered a sequence of cascading failures that brought down a datacenter and the customer machines hosted there; a corrupted routing table at GoDaddy disconnected their domain name servers for a day and caused a widespread outage; and a network connectivity issue at United Airlines took down their reservation system, leading to thousands of flight cancellations and a "ground stop" at their San Francisco hub.

One way to make networks more reliable would be to develop tools for checking important network invariants automatically. These tools would allow administrators to answer questions such as: "does this configuration provide connectivity to every host

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

PLDI'13, June 16–19, 2013, Seattle, WA, USA. Copyright © 2013 ACM 978-1-4503-2014-6/13/06...\$15.00

in the network?" or "does this configuration correctly enforce the access control policy?" or "does this configuration have a forwarding loop?" or "does this configuration properly isolate trusted and untrusted traffic?" Unfortunately, until recently, building such tools has been effectively impossible due to the complexity of today's networks. A typical enterprise or datacenter network contains thousands of heterogeneous devices, from routers and switches, to web caches and load balancers, to monitoring middleboxes and firewalls. Moreover, each device executes a stack of complex protocols and is configured through a proprietary and idiosyncratic interface. To reason formally about such a network, an administrator (or tool) must reason about the proprietary programs running on each distributed device, as well as the asynchronous interactions between them. Although formal models of traditional networks exist, they have either been too complex to allow effective reasoning, or too abstract to be useful. Overall, the incidental complexity of networks has made reasoning about their behavior practically infeasible.

Fortunately, recent years have seen growing interest in a new kind of network architecture that could provide a foundation for network reasoning. In a *software-defined network* (SDN), a program on a logically-centralized *controller machine* defines the overall policy for the network, and a collection of *programmable switches* implement the policy using efficient packet-processing hardware. The controller and switches communicate via an open and standard interface. By carefully installing packet-processing rules in the hardware tables provided on switches, the controller can effectively manage the behavior of the entire network.

Compared to traditional networks, SDNs have two important simplifications that make them amenable to formal reasoning. First, they relocate control from distributed algorithms running on individual devices to a single program running on the controller. Second, they eliminate the heterogeneous devices used in traditional networks—switches, routers, load balancers, firewalls, etc.—and replace them with stock programmable switches that provide a standard set of features. Together, this means that the behavior of the network is determined solely by the sequence of configuration instructions issued by the controller. To verify that the network has some property, an administrator (or tool) simply has to reason about the states of the switches as they process instructions.

In the networking community, there is burgeoning interest in tools for checking network-wide properties automatically. Systems such as FlowChecker [1], Header Space Analysis [12], Anteater [17], VeriFlow [13], and others, work by generating a logical representation of switch configurations and using an automatic solver to check properties of those configurations. The configurations are obtained by "scraping" state off of the switches or inspecting the instructions issued by an SDN controller at run-time.

These tools represent a good first step toward making networks more reliable, but they have two important limitations. First, they are based on ad hoc foundations. Although SDN platforms such as OpenFlow [21] have precise (if informal) specifications, the tools make simplifying assumptions that are routinely violated by real

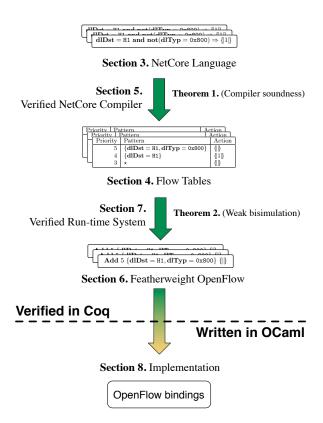


Figure 1: System architecture.

hardware switches. For example, some tools assume that switches will process instructions emitted by the controller in sequence, even though actual switches often reorder messages. This means that properties they consider verified do not always actually hold. Second, the tools are expensive to run and do not scale well. For example, most tools take several minutes to run, even in small to medium-sized networks (VeriFlow [13] is a notable exception). This is too slow to be used in large dynamic networks where configurations change on the order of seconds. Overall, although these tools are useful for finding bugs, they cannot provide the rigorous guarantees that networks, as critical infrastructure, require.

Our approach. This paper presents a different approach. Rather than building tools to find bugs in SDN controllers at run-time, we develop a verified SDN controller in the Coq proof assistant and prove it correct against a formal specification and a detailed operational model of SDN. With our controller, programmers specify the behavior of the network using the NetCore programming language [22], which abstracts away from the details of the underlying switch hardware and distributed system, and allows programmers to reason in terms of simple hop-by-hop packet-processing steps. The NetCore compiler and run-time system translates programs written in this language down to low-level packet-processing rules. Because its behavior is verified in Coq, we establish the correctness of our controller once and for all, obviating the need for run-time or post hoc verification as in most current tools.

Architecturally, our system is organized as a verified software stack that translates through the following levels of abstraction:

NetCore. The highest level of abstraction is the NetCore language, proposed in prior work by Monsanto et al. [22]. NetCore is a declarative language that allows programmers to describe what network behavior they want, without specifying

how it should be implemented. It offers a collection of intuitive constructs for matching, filtering, and transforming packets, as well as natural logical operators for combining smaller programs into bigger ones such as union and domain restriction. Although NetCore programs are ultimately executed in a distributed system—the network—they have a simple semantics that models their behavior as functions from packets to packets.

- Flow tables. The intermediate level of abstraction is *flow tables*, a representation that sits between NetCore programs and switch-level configurations. There are two main differences between NetCore programs and flow tables. First, NetCore programs describe the forwarding behavior of a whole network, while flow tables describe the behavior of a single switch. Second, flow tables process packets using a linear scan through a list of prioritized rules. Hence, to translate operators such as union and negation, the NetCore compiler must generate a sequence of rules that encodes the same semantics. However, because flow table matching uses a lower-level packet representation (as nested frames of Ethernet, IP, TCP, etc. packets), flow tables must satisfy a well-formedness condition to rule out invalid patterns that are inconsistent with this representation.
- Featherweight OpenFlow. The lowest level of abstraction is Featherweight OpenFlow, a new foundational model we have designed that captures the essential features of SDNs. Featherweight OpenFlow models switches, the controller, the network topology, as well as their internal transitions and interactions in a small-step operational semantics. This semantics is nondeterministic, modeling the asynchrony inherent in networks. To implement a flow table in a Featherweight OpenFlow network, the controller instructs switches to install or uninstall rules as appropriate while dealing with two important issues: First, switches process instructions concurrently with packets flowing through the network, so it must ensure that at all times the rules installed on switches are consistent with the flow table. Second, switches are allowed to buffer instructions and apply them in any order, so it must ensure that the behavior is correct no matter how instructions are reordered through careful use of synchronization primitives.

Figure 1 depicts the architecture of our system, and provides an outline for this paper. Overall, our main contributions are as follows:

- We present the first machine-verified SDN controller, which gives network programmers robust static guarantees backed by machine-checked proofs against a foundational model.
- We develop Featherweight OpenFlow, the first formal model of OpenFlow. It includes all sources of asynchrony and nondeterminism mentioned in the informal OpenFlow specification, as well as a precise model of switch flow table semantics.
- We formalize NetCore, flow tables, and Featherweight Open-Flow in Coq, and develop machine-checked proofs of correctness for the translations between them.
- We present our prototype implementation, obtained by extracting our Coq development to OCaml, and present experimental results comparing the performance of our system against unverified controllers on simple benchmarks.

Besides their use in our system, we hope that the abstractions and theorems presented in this paper will be useful to others. Flow tables are a canonical representation of switch state that appear in many other systems. Likewise, Featherweight OpenFlow is a comprehensive model that captures the essential forwarding behavior of SDNs in a minimal core calculus. Our design and Coq formalization of flow tables and Featherweight OpenFlow provide

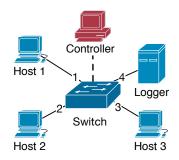


Figure 2: Example network topology.

a starting point for developing extensions that model additional SDN features and a foundation for other verified systems.

2. Overview

To motivate the need for verified SDN controllers, consider the simple network depicted in Fig. 2. It consists of a single switch connected to four hosts: three clients and a middlebox that monitors HTTP requests. The ports on the switch are numbered 1 to 4; the clients are connected to ports 1 to 3 and the middlebox to port 4. In addition, the switch has a dedicated link to the controller that is not considered part of the data network.

Now imagine we want to build an SDN controller that implements the following high-level network policy: block SSH traffic, log HTTP requests, and allow clients to send non-SSH traffic to each other. It is straightforward to formalize this policy as a packet-processing function that maps input packets to (possibly several) output packets: the function drops SSH packets, forwards HTTP packets both to their destination and to the middlebox, and forwards all other packets to their destination alone.

To implement this function in an SDN, however, we would need to specify several additional low-level details, since switches cannot implement general packet-processing functions directly. First, the controller would need to encode the function as a *flow table*—a set of prioritized forwarding rules. Second, it would need to send the switch a series of control messages to add individual entries from the flow table, incrementally building up the complete table.

More concretely, the controller could first send a message instructing the switch to add a flow table entry that blocks SSH traffic:

Add
$$10 \{ tpDst = 22 \} \{ | \}$$

Here 10 is a priority number, $\{\mathbf{tpDst} = 22\}$ is a pattern that matches SSH traffic (TCP port 22), and $\{\|\}$ is an empty multiset of ports, which drops packets, as intended. Next, the controller could add an entry to process HTTP requests:

Add 9 {**dlDst** =
$$H1$$
, **tpDst** = 80} { $[1, 4]$ }

Note that this rule duplicates HTTP (TCP port 80) packets, sending them to the monitor and to their destination. Finally, the controller could add an entry to forward other packets to their destination:

Add 1 {**dlDst** =
$$H1$$
} {|1|}

Note that this rule does not apply to SSH and HTTP traffic, since those packets are handled by the higher-priority rules.

After these control messages have been sent, it would be natural to expect that the network correctly implements the packetprocessing function described above. But the situation is actually Packet pk ::=**Eth** dlSrc dlDst dlTyp nwPkNetwork layer nwPk ::=**IP** nwSrc nwDst nwProto tpPk

Transport layer tpPk ::= TCP tpSrc tpDst payload| Unknown payload

Figure 3: Logical packet structure.

more complicated: switches have substantial latitude in how they process messages from the controller, and packets may arrive at any time during processing. Establishing that the network correctly implements this function—in particular, that it blocks SSH traffic and logs HTTP traffic—requires additional reasoning.

Controller-switch consistency. Switches process packets and control messages concurrently. In our example, the switch may receive an HTTP request before the flow table entry that handles HTTP packets arrives. In this case, the switch will send the packet to the controller for further processing. Since the controller is a general-purpose machine, it can implement the packet-processing function directly, apply it to the incoming packet, and send the results back to the switch. However, this means that SDN controllers typically have two different implementations of the function: one residing at the controller and another on the switches. A key property we verify is that these two implementations are consistent.

Message reordering. SDN switches may process control messages in any order, and many switches do, to maximize performance. But unrestricted reordering can cause implementations to violate their intended specifications. For example, if the rule to drop SSH traffic is installed after the final, low-priority rule that forwards all traffic, then SSH traffic will temporarily be forwarded by the low-priority rule, breaking the intended security policy. To ensure that such reorderings do not occur, a controller must carefully insert barrier messages, which force the switch to process all outstanding messages. A key property we verify is that controllers use barriers correctly (several unverified controllers ignore this issue).

Natural patterns. Another complication is that the patterns presented earlier in this section, such as $\{tpDst = 22\}$, are actually invalid. To match SSH traffic, it is not enough to simply state that the destination port must be 22. The pattern must also specify that the Ethernet frame type must be IP, and the transport protocol must be TCP. Without these additional constraints, switches will interpret the pattern as a wildcard that matches all packets. Several earlier controller platforms did not properly account for this behavior, and had bugs as a result. We develop a semantics for patterns and identify a class of natural patterns that are closed under the algebraic operations used by our compiler and flow table optimizer.

Roadmap. The rest of this paper develops techniques for establishing that a given packet-processing function is implemented correctly by an OpenFlow network. More specifically, we tackle the problem of verifying high-level programming abstractions, using NetCore [22] as a concrete instance of a high-level network language. The next section presents NetCore in detail. The following sections describe general and reusable techniques for establishing the correctness of SDN controllers, including NetCore.

3. NetCore

This section presents the highest layer of our verified stack: the NetCore language. A NetCore program specifies how the switches process packets at each hop through the network. More formally, a program denotes a total function from port-packet pairs to multisets of port-packet pairs. The syntax and semantics of a core NetCore

 $^{^{1}}$ The controller would actually need to create rules for each client. To save space, we have only given the rules for HI here.

```
Switch ID
                sw \in \mathbb{N}
Port ID
                pt \in \mathbb{N}
                h ::= dlSrc \mid dlDst
Headers
                                                      MAC address
                                                      Ethernet frame type
                         dlTyp
                         nwSrc | nwDst
                                                      IP address
                         nwProto
                                                      IP protocol code
                         tpSrc | tpDst
                                                      transport port
Predicate pr := \star
                                                      wildcard
                        h = n
                                                      match header
                                                      match switch
                                                      predicate negation
 \begin{array}{c|c} & pr_1 \text{ and } pr_2 \\ \text{Program} & pg ::= pr \Rightarrow \{|pt_1 \cdots pt_n|\} \end{array} 
                                                      predicate conjunction
                                                      basic program
                                                      program union
                      | pg_1 \uplus pg_2
                         restrict pq by pr
                                                      program restriction
 [pr] sw pt pk
   [\star] sw pt pk = true
   \|\mathbf{d}\|\mathbf{Src} = n\| sw pt (Eth dlSrc _ _ _ ) = dlSrc = n
  \llbracket \mathbf{nwSrc} = n \rrbracket \ sw \ pt \ (\mathbf{Eth} \_\_\_ (\mathbf{IP} \ nwSrc \_\_\_)) = nwSrc = n
  [\![\mathbf{nwSrc} = n]\!] sw pt (Eth _ _ _ (Unknown _)) = false
  \llbracket \mathbf{at} \ sw' \rrbracket \ sw \ pt \ pk = sw = sw'
  [pg] sw pt pk = \{ [(pt_1, pk_1) \cdots (pt_n, pk_n)] \}
  \llbracket pg_1 \uplus pg_2 \rrbracket \ sw \ pt \ pk =
     \llbracket \bar{p}g_1 \rrbracket \ \bar{sw} \ pt \ p\hat{k} \ \dot{\oplus} \ \llbracket pg_2 \rrbracket \ sw \ pt \ pk
  [restrict pg by pr] sw pt pk =
```

Figure 4: NetCore syntax and semantics (extracts).

 $\{|(pt',pk')| (pt',pk') \in \llbracket pg \rrbracket \ sw \ pt \ pk \land \llbracket pr \rrbracket \ sw \ pt \ pk \}$

fragment are shown in Fig. 4. To save space, we have elided several header fields and operators not used in this paper.

We can build a NetCore program that implements the example from the previous section by composing several smaller NetCore program fragments. The first fragment forwards traffic to *H1*:

$$pg_1 \triangleq \mathbf{dlDst} = HI \Rightarrow \{|1|\}$$

This basic program consists of a predicate pr and a multiset of actions $\{pt_1 \cdots pt_n\}$. The predicate denotes a set of port-packet pairs, and the actions denote the ports (if any) where those packets should be forwarded on the next hop. In this instance, the predicate denotes the set of all packets whose Ethernet destination (dlDst) address has the specified value, and the actions denote a transformation that forwards matching packets to port 1. Note that we represent packets as nested sequences of frames (Ethernet, IP, TCP, etc.) as shown in Fig. 3. NetCore provides predicates for matching on well-known header fields as well as logical operators such as and and or, unlike hardware switches, which only provide prioritized sets of rules.

The next two fragments are similar to pg_1 , but forward traffic to H2 and H3 instead of H1:

$$\begin{array}{l} pg_2 \stackrel{\triangle}{=} \mathbf{dlDst} = H2 \Rightarrow \{\![2]\!\} \\ pg_3 \stackrel{\triangle}{=} \mathbf{dlDst} = H3 \Rightarrow \{\![3]\!\} \end{array}$$

Using the union operator, we can combine these programs into a single program that implements forwarding between all clients:

$$pg_{\mathit{fwd}} \triangleq pg_1 \uplus pg_2 \uplus pg_3$$

Semantically, the \uplus operator produces the (multiset) union of the results produced by each sub-program. Using the union operator again, we can extend this program to one that also forwards HTTP requests to the middlebox:

$$pg_{fwd} \uplus \mathbf{tpDst} = 80 \Rightarrow \{ |4| \}$$

Note that this program duplicates packets sent to port 80, forwarding to their destination and also to the logging machine. Finally, we can add the security policy using the **restrict by** operator, which restricts a program by a predicate:

restrict
$$(pg_{fwd} \uplus \mathbf{tpDst} = 80 \Rightarrow \{4\}) \mathbf{by} (\mathbf{not} \mathbf{tpDst} = 22)$$

This program is similar the previous one, but drops SSH traffic.

The advantages of using a declarative language such as NetCore should be clear: it provides abstractions that make it easy to establish network-wide properties through compositional reasoning. For example, simply by inspecting the final program and using the denotational semantics (Fig. 4), we can easily verify that the network blocks SSH traffic, forwards HTTP traffic to the middlebox, and provides pair-wise connectivity between the clients. In particular, even though a controller, switches, flow tables, forwarding rules, are all involved in implementing this program, we do not have to reason about them! This is in contrast to lower-level controller platforms. which require programmers to explicitly construct switch-level forwarding rules, issue messages to install those rules on switches, and reason about the asynchronous interactions between switches and controller. Of course, the complexity of the underlying system is not eliminated, but relocated from the programmer to the language implementers. This is an efficient tradeoff: functionality common to many programs can be implemented just once, proved correct, and reused broadly.

4. Flow Tables

The first step toward executing a NetCore program in an SDN involves compiling it to a prioritized set of forwarding rules—a *flow table*. Flow tables are an intermediate representation that play a similar role in NetCore to register transfer language (RTL) in traditional compilers. Flow tables are more primitive than NetCore programs because they lack the logical structure induced by NetCore operators such as union, intersection, negation, and restriction. Also, the patterns used to match packets in flow tables are more restrictive than NetCore predicates. And unlike NetCore programs, which denote total functions, flow tables are partial: switches redirect unmatched packets to the controller.

As defined in Fig. 5, a *flow table* consists of a multiset of rules (n, pat, pts) where n is an integer priority, pat is a pattern, and pts is a multiset of ports. A *pattern* is a record that associates each header field to either an integer constant n or the special *wildcard* value \star . When writing flow tables, we often elide headers set to \star in patterns as well as priorities when they are clear from context.

Pattern semantics. The semantics of patterns is given by the function pk#pat, as defined in Fig. 5. This turns out to be subtly complicated, due to the representation of packets as sequences of nested frames—a pattern contains a (possibly wildcarded) field for every header field, but not all packets contain every header field. Some fields only exist in specific frame types (**dlTyp**) or protocols (**nwProto**). For example, only IP packets (**dlTyp** = 0x800) have IP source and destination addresses. Likewise, TCP (**nwProto** = 6) and UDP (**nwProto** = 17) packets have source and destination ports, but ICMP (**nwProto** = 1) packets do not.

```
Wildcard
                    w ::= n \mid \star
Pattern
                    pat ::= \{ \mathbf{dlSrc} = w, \mathbf{dlDst} = w, \mathbf{dlTyp} = w, 
                                   \mathbf{nwSrc} = w, \mathbf{nwDst} = w, \mathbf{nwProto} = w,
                                   tpSrc = w, tpDst = w
Flow table FT \in \{|n \times pat \times \{|pt|\}|\}
\llbracket FT \rrbracket pt pk \rightsquigarrow \{ pt_1 \cdots pt_n \} \times \{ pk_1 \cdots pk_m \}
                       \exists (n, pat, \{|pt_1 \cdots pt_n|\}) \in FT.
                   pk\#pat = \mathbf{frue}
\forall (n', pat', pts') \in FT. \ n' > n \Rightarrow
pk\#pat' = \mathbf{false}
[FT] \ pt \ pk \leadsto (\{(pt_1) \cdots (pt_n)\}, \{\}\})  (MATCHED)
         \frac{\forall (n, pat, pts) \in FT \qquad pk\#pat = \textbf{false}}{\|FT\| \ pt \ pk \leadsto (\{\|\}, \{|(pt, pk)\})} \ (\text{Unmatched})
pk\#pat
               (Eth dlSrc dlDst dlTyp nwPk)#pat =
                   dlSrc \sqsubseteq pat.dlSrc \land dlDst \sqsubseteq pat.dlDst \land
                   dlTyp \sqsubseteq pat.\mathbf{dlTyp} \wedge
                   (pat.\mathbf{dlTyp} = 0x800 \Rightarrow nwPk\#_{nw}pat)
 nwPk\#_{nw}pat
             (IP nwSrc\ nwDst\ nwProto\ tpPk)\#_{nw}pat =
                nwSrc \sqsubseteq pat.\mathbf{nwSrc} \land nwDst \sqsubseteq pat.\mathbf{nwDst} \land
                nwProto \sqsubseteq pat.\mathbf{nwProto} \land
                 (pat.\mathbf{nwProto} = 6 \Rightarrow tpPk\#_{tp}pat)
             (Unknown payload)\#_{nw} pat = true
tpPk\#_{tp}pat
                 (TCP tpSrc\ tpDst\ payload)#_{tp} pat =
                     tpSrc \sqsubseteq pat.\mathbf{tpSrc} \land tpDst \sqsubseteq pat.\mathbf{tpDst}
                 Unknown payload \#_{tp} pat = true
n \sqsubseteq w
                       m \sqsubseteq n = m = n n \sqsubseteq \star = true
```

Figure 5: Flow table syntax and semantics.

To match on a given field, a pattern must specify values for all other fields it depends on. For example, to match on IP addresses, the pattern must also specify that the Ethernet frame type is IP:

$$\{dlTyp = 0x800, nwSrc = 10.0.0.1\}$$

If the frame type is elided, the value of the dependent header is silently ignored and the pattern is equivalent to a wildcard:

$$\{nwSrc = 10.0.0.1\} \equiv \{\}$$

In effect, patterns not only match packets, but also determine how they are parsed. This behavior, which was ambiguous in early versions of the OpenFlow specification (and later fixed) has lead to real bugs in existing controllers (Section 5). Although unintuitive for programmers, this behavior is completely consistent with how packet processing is implemented in modern switch hardware.

Flow table semantics. The semantics of flow tables is given by the relation $[\![\cdot]\!]$. The relation has two cases: one for matched packets and another for unmatched packets. Each flow table entry is a tuple containing a priority n, pattern pat, and a multiset of ports

```
\mathcal{P}: sw \times pr \rightarrow [(pat, \mathbf{bool})]
 \begin{array}{ll} \mathcal{P}(sw, \mathbf{dlSrc} = n) &= \left[ (\{\mathbf{dlSrc} = n\}, \mathbf{true}) \right] \\ \mathcal{P}(sw, \mathbf{nwSrc} = n) &= \left[ (\{\mathbf{dlTyp} = 0 \mathrm{x800}, \mathbf{nwSrc} = n\}, \mathbf{true}) \right] \end{array}
                                                     = [(\star, true)]
  \mathcal{P}(sw, \mathbf{at}\ sw)
                                                    = [(\star, \mathbf{false})] \quad \text{where } sw \neq sw'
= [(pat_1, \neg b_1) \cdots (pat_n, \neg b_n), (\star, \mathbf{false})]
  \mathcal{P}(sw, \mathbf{at} \ sw')
  \mathcal{P}(sw, \mathbf{not}\ pr)
           where [(pat_1, b_1) \cdots (pat_n, b_n)] = \mathcal{P}(sw, pr)
  \mathcal{P}(sw, pr \text{ and } pr') =
           \begin{aligned} &[(pat_1 \cap pat_1', b_1 \wedge b_1') \cdots (pat_m \cap pat_n', b_m \wedge b_n')] \\ &\text{where } [(pat_1, b_1) \cdots (pat_m, b_m)] = \mathcal{P}(sw, pr) \\ &\text{where } [(pat_1', b_1') \cdots (pat_n', b_n')] = \mathcal{P}(sw, pr') \end{aligned} 
 C: sw \times pg \rightarrow [(pat, pt)]
C(sw, pr \Rightarrow pt) = [(pat_1, pt_1) \cdots (pat_n, pt_n), (\star, \{\{\}\})]
       where [(pat_1, b_1), \cdots, (pat_n, b_n)] = \mathcal{P}(sw, pr)
      where pt_i = pt if b_i = \mathbf{true}
      where pt_i = \{|\} if b_i = false
 \begin{array}{c} \mathcal{C}(sw,pg \uplus pg') = \\ [(pat_1 \cap pat'_1,pt_1 \uplus pt'_1),\cdots,(pat_m \cap pat'_n,pt_m \uplus pt'_n)] + \\ \end{array} 
      \begin{aligned} &[(pat_1,pt_1)\cdots(pat_n,pt_m)] + \\ &[(pat_1',pt_1')\cdots(pat_n',pt_n')] + \\ &[(pat_1',pt_1')\cdots(pat_n',pt_n')] = \mathcal{P}(sw,pg) \\ &\text{where } [(pat_1',pt_1')\cdots(pat_n',pt_n')] = \mathcal{P}(sw,pg') \end{aligned}
```

Figure 6: NetCore compilation.

 $\{|pt_1\cdots pt_n|\}$. Given a packet and its input port, the semantics forwards the packet to all ports in the multiset associated with the highest-priority matching rule in the table. Otherwise, if no matching rule exists, it diverts the packet to the controller. In the formal semantics, the first component of the result pair represents forwarded packets while the second component represents diverted packets. Note that flow table matching is non-deterministic if there are multiple matching entries with the same priority. This has serious implications for a compiler—e.g., naively combining flow tables with overlapping priorities could produce incorrect results. In the NetCore compiler, we avoid this issue by always working with unambiguous and total flow tables.

5. Verified NetCore Compiler

With the syntax and semantics of NetCore and flow tables in place, we now present a verified compiler for NetCore. The compiler takes programs as input and generates a set of flow tables as output, one for every switch. The compilation algorithm is based on previous work [22], but we have verified its implementation in Coq. While building the compiler, we found two serious bugs in the original algorithm related to the handling of (unnatural) patterns in the compiler and flow table optimizer.

The compilation function \mathcal{C} , defined in Fig. 6, generates a flow table for a given switch sw. It uses the auxiliary function \mathcal{P} to compile predicates. The compiler produces a list of pattern-action pairs, but priority numbers are implicit: the pair at the head has highest priority and each successive pair has lower priority.

Because NetCore programs denote total functions, packets not explicitly matched by any predicate are dropped. In contrast, flow tables divert unmatched packets to the controller. The compiler resolves this discrepancy by adding a catch-all rule that drops unmatched packets. For instance:

```
C(sw, dlSrc = H1 \Rightarrow \{|5|\}) = \{|(2, \{dlSrc = H1\}, \{|5|\}), (1, \star, \{|\})\}\}
```

The key operator used by the compiler constructs the cross-product of the flow tables provided as input. This operator can be used to compute intersections and unions of flow tables. Note that implementing union in the "obvious" way—by concatenating flow tables—would be wrong. The cross-product operator performs an element-wise intersection of the input flow tables and then merges their actions. To compile a union, we first use cross-product to build a flow table that represents the intersection, and then concatenate the flow tables for the sub-policies at lower priority. For example, the following NetCore program,

$$dlSrc = H1 \Rightarrow \{|5|\} \uplus dlDst = H2 \Rightarrow \{|10|\}$$

compiles to a flow table:

Priority	Pattern	Action
4	$\{\mathbf{dlSrc} = H1, \mathbf{dlDst} = H2\}$	$\{ [5, 10] \}$
3	$\{\mathbf{dlSrc} = H1\}$	$\{ 5 \}$
2	$\{\mathbf{dlDst} = H2\}$	$\{ 10 \}$
1	*	{ }

The first rule matches all packets that match both sub-programs, while the second and third rules match packets only matched by the left and the right programs respectively. The final rule drops all other packets. The compilation of other predicates uses similar manipulations on flow tables.

We have built a large library of flow table manipulation operators in Coq, along with several lemmas that state useful algebraic properties about these operators. With this library, proving the correctness theorem for the NetCore compiler is simple—only about 200 lines of code in Coq.

Theorem 1 (NetCore Compiler Soundness). For all NetCore programs pg, switches sw, ports pt, and packets pk we have ||C(sw, pg)|| pt pk = ||pg|| sw pt pk.

Intuitively, this theorem states that a flow table compiled from a NetCore program for a switch sw, has the same behavior as the NetCore program evaluated on packets at sw.

Compiler bugs. In the course of our work, we discovered that several unverified compilers from high-level network programming languages to flow tables suffer from bugs due to subtle pattern semantics. Section 4 described inter-field dependencies in patterns. For example, to match packets from IP address 10.0.0.1, we write

$$\{nwSrc = 10.0.0.1, dlTyp = 0x800\}$$

and if we omit the **dlTyp** field, the IP address is silently ignored. This unintuitive behavior has led to bugs in PANE [6] and Nettle [27] as well as an unverified version of NetCore [22]. To illustrate, consider the following program:

$$nwSrc = 10.0.0.1 \Rightarrow \{|5|\}$$

In NetCore, this program matches all IP packets from 10.0.0.1 and forwards them out port 5. But the original NetCore compiler produced the following flow table for this program:

Priority	Pattern	Action
2	$\{$ nwSrc = 10.0.0.1 $\}$	{ 5 }
1	*	{ }

Because the first pattern is equivalent to the all-wildcard pattern, this flow table sends *all* traffic out port 5. Both PANE and Nettle have similar bugs. Nettle has a special case to handle patterns with IP addresses that do not specify **dlTyp** = 0x800, but it does not correctly handle patterns that specify a transport port number but not the **nwProto** field. PANE suffers from the same bug. Even worse, these invalid patterns lead to further bugs when flow tables are combined and optimized by the compiler.

Natural patterns. The verified NetCore compiler does not suffer from the bug above. In our formal development, we require that

all patterns manipulated by the compiler be what we call *natural* patterns. A natural pattern has the property that if the pattern specifies the value of a field, then all of that field's dependencies are also met. This rules out patterns such as $\{nwSrc = 10.0.0.1\}$, which omits the Ethernet frame type necessary to parse the IP address. Natural patterns are easy to define using dependent types in Coq. Moreover, we can calculate the cross-product of two natural patterns by intersecting fields point-wise. Hence, it is easy to prove that natural patterns are closed under intersection.

Lemma 1. If pat_1 and pat_2 are natural patterns, then $pat_1 \cap pat_2$ is also a natural pattern.

Another important property is that all patterns can be expressed as some equivalent natural pattern (where patterns are equivalent if they denote the same set of packets). This property tells us that we do not lose expressiveness by restricting to natural patterns.

Lemma 2. If pat is an arbitrary pattern, then there exists a natural pattern pat', such that pat \equiv pat'.

These lemmas are used extensively in the proofs of correctness for our compiler and flow table optimizer.

Flow table optimizer. The basic NetCore compilation algorithm described so far generates flow tables that correctly implement the semantics of the input program. But many flow tables have redundant entries that could be safely removed. For example, a naive compiler might translate the program $(\star \Rightarrow \{5\})$ to the flow table $\{(2,\star,\{5\}),(1,\star,\{\})\}\}$, which is equivalent to $\{(2,\star,\{5\})\}$. Worse, because the compilation rule for union uses a cross-product operator to combine the flow tables computed for sub-programs, the output can be exponentially larger than the input. Without an optimizer, such a naive compiler is essentially useless—e.g., we built an unoptimized implementation of the algorithm in Fig. 6 and found that it ran out of memory when compiling a program consisting of just 9 operators!

Our compiler is parameterized on a function $\mathcal{O}: FT \to FT$, that it invokes at each recursive call. Because even simple policies can see a combinatorial explosion during compilation, this inline reduction is necessary. We stipulate that \mathcal{O} must produce equivalent flow tables: $\llbracket \mathcal{O}(FT) \rrbracket = \llbracket FT \rrbracket$.

We have built an optimizer that eliminates low-priority entries whose patterns are fully subsumed by higher-priority rules and proved that it satisfies the above condition in Coq. Although this optimization is quite simple, it is effective in practice. In addition, earlier attempts to implement this optimization in NetCore had a bug that incorrectly identified certain rules as overlapping which we did not discover until developing this proof. The PANE optimizer also had a bug—it assumed that combining identical actions is always idempotent. Both of these bugs led to incorrect behavior.

6. Featherweight OpenFlow

The next step towards executing NetCore programs is a controller that configures the switches in the network. To prove that such a controller is correct, we need a model of the network. Unfortunately, the OpenFlow 1.0 specification, consisting of 42 pages of informal prose and C definitions, is not amenable to rigorous proof.

This section presents Featherweight OpenFlow, a detailed operational model that captures the essential features of OpenFlow networks, and yet still fits on a single page. The model elides a number of details such as error codes, counters, packet modification, and advanced configuration options such as the ability to enable and disable ports. But it does include all of the features related to how packets are forwarded and how flow tables are modified. Many existing SDN bug-finding and property-checking tools are based on

```
Ports on switch
                                S ::= \mathbb{S}(sw, pts, FT, in_p, out_p, in_m, out_m)
                                                                                                                                                              \in \{pt\}
     Switch
                                                                                                                                                 in_p, out_p \in \{ (pt, pk) \}

in_m \in \{ SM \}
     Controller
                                C ::= \mathbb{C}(\sigma, f_{in}, f_{out})
                                                                                                           Input/output buffers
     Link
                                L ::= \mathbb{L}((sw_{src}, pt_{src}), pks, (sw_{dst}, pt_{dst}))
                                                                                                           Messages from controller
     Link to Controller
                               M := \mathbb{M}(sw, SMS, CMS)
                                                                                                           Messages to controller
                                                                                                                                                              \in \{|CM|\}
                                                                                                                               Switch Components
                                           Devices
           Controller state
                                                                                                         Message queue from controller SMS \in [SM_1 \cdots SM_n]
           Controller input relation f_{in} \in sw \times CM \times \sigma \leadsto \sigma
Controller output relation f_{out} \in \sigma \leadsto sw \times SM \times \sigma
                                                                                                         Message queue to controller CMS \in [CM_1 \cdots CM_n]
                                                                                                                                  Controller Link
                               Controller Components
                                            From controller
                                                                    SM ::= FlowMod \delta \mid PktOut pt \mid pk \mid BarrierRequest n
                                                                     CM ::= \mathbf{PktIn} \ pt \ pk \mid \mathbf{BarrierReply} \ n
                                            To controller
                                            Table update
                                                                    \delta ::= Add n pat act | Del pat
                                                                           Abstract OpenFlow Protocol
                                                             [FT](pt, pk) \leadsto (\{pt'_1 \cdots pt'_n\}, \{pk'_1 \cdots pk'_m\})
                                                                                                                                                                                  (FWD)
                                 \mathbb{S}(sw, pts, FT, \{|(pt, pk)|\} \uplus in_p, out_p, in_m, out_m)
                                 S(sw, pts, FT, in_p, \{(pt'_1, pk) \cdots (pt'_n, pk)\} \uplus out_p, in_m, \{\{PktIn \ pt \ pk'_1 \cdots PktIn \ pt \ pk'_m\} \uplus out_m)
                                                                                                                                                                         (WIRE-SEND)
                                             \mathbb{S}(sw, pts, FT, in_p, \{ (pt, pk) \} \oplus out_p, in_m, out_m) \mid \mathbb{L}((sw, pt), pks, (sw', pt'))
                                          \mathbb{S}(sw, pts, FT, in_p, out_p, in_m, out_m) \mid \mathbb{L}((sw, pt), [pk] + pks, (sw', pt'))
                                                                                                                                                                        (WIRE-RECV)
                                             \mathbb{L}((sw',pt'),pks+|pk|,(sw,pt)) \mid \mathbb{S}(sw,pts,FT,in_p,out_p,in_m,out_m)
                                    \longrightarrow \mathbb{L}((sw',pt'),pks,(sw,pt)) \mid \mathbb{S}(sw,pts,FT,\{[(pt,pk)]\} \uplus in_p,out_p,in_m,out_m)
\mathbb{S}(\overline{sw,pts,FT,in_p,out_p,\{\textbf{FlowMod Add}\ m\ pat\ act\}}\ \uplus\ in_m,out_m)\longrightarrow \mathbb{S}(sw,pts,FT\ \uplus\ \{[(m,pat,act)]\},in_p,out_p,in_m,out_m))))
                                                FT_{rem} = \{ | (n', pat', act') \mid (n', pat', act') \in FT \text{ and } pat \neq pat' | \}
                                                                                                                                                                                   (DEL)
                       \mathbb{S}(sw, pts, FT, in_p, out_p, \{\{FlowMod Del pat\}\} \uplus in_m, out_m) \longrightarrow \mathbb{S}(sw, pts, FT_{rem}, in_p, out_p, in_m, out_m)
                                                                                         pt \in pts
                 \overline{\mathbb{S}(sw,pts,FT,in_p,out_p,\{\![\mathbf{PktOut}\;pt\;pk]\!] \uplus in_m,out_m)} \longrightarrow \mathbb{S}(sw,pts,FT,in_p,\{\![(pt,pk)]\!] \uplus out_p,in_m,out_m)} \ (\mathsf{PKTOUT})
                                                                              f_{out}(\sigma) \leadsto (sw, SM, \sigma')
                                                                                                                                                                         (CTRL-SEND)
                                 \overline{\mathbb{C}(\sigma, f_{in}, f_{out}) \mid \mathbb{M}(sw, SMS, CMS) \longrightarrow \mathbb{C}(\sigma', f_{in}, f_{out}) \mid \mathbb{M}(sw, [SM] + SMS, CMS)}
                                                                               f_{in}(sw, \sigma, CM) \leadsto \sigma'
                                                                                                                                                                        (CTRL-RECV)
                                 \overline{\mathbb{C}(\sigma, f_{in}, f_{out}) \mid \mathbb{M}(sw, SMS, CMS + [CM]) \longrightarrow \mathbb{C}(\sigma', f_{in}, f_{out}) \mid \mathbb{M}(sw, SMS, CMS)}
                                                                             SM \neq \mathbf{BarrierRequest} \ n
                                                                                                                                                            (SWITCH-RECV-CTRL)
                                                    \mathbb{M}(sw, SMS + [SM], CMS) \mid \mathbb{S}(sw, pts, FT, in_p, out_p, in_m, out_m)
                                                   \mathbb{M}(sw, SMS, CMS) \mid \mathbb{S}(sw, pts, FT, in_p, out_p, \{SM\} \uplus in_m, out_m)
                                                                                                                                                       (SWITCH-RECV-BARRIER)
                                   \mathbb{M}(sw, SMS + | \mathbf{BarrierRequest} \ n | , CMS) | \mathbb{S}(sw, pts, FT, in_p, out_p, \{ | \}, out_m )
                                  \mathbb{M}(sw, SMS, CMS) \mid \mathbb{S}(sw, pts, FT, in_p, out_p, \{\}\}, \{\{\}\}\}
                                                                                                                                                            (SWITCH-SEND-CTRL)
                                                    \mathbb{S}(sw, pts, FT, in_p, out_p, in_m, \{CM\} \uplus out_m) \mid \mathbb{M}(sw, SMS, CMS)
                                           \longrightarrow \mathbb{S}(sw, pts, FT, in_p, out_p, in_m, out_m) \mid \mathbb{M}(sw, SMS, [CM] + CMS)
                                                                          \frac{Sys_1 \longrightarrow Sys_1'}{Sys_1 \mid Sys_2 \longrightarrow Sys_1' \mid Sys_2}
                                                                                                                                                                      (CONGRUENCE)
```

Figure 7: Featherweight OpenFlow syntax and semantics.

similar (informal) models [3, 12, 13]. We believe Featherweight OpenFlow could also serve as a foundation for these tools.

6.1 OpenFlow Semantics

Initially, every switch has an empty flow table that diverts all packets to the controller. Using **FlowMod** messages, the controller can insert new table entries to have the switch process packets itself. A non-trivial program may compile to several thousand flow table entries, but **FlowMod** messages only add a single entry at a time. In general, many **FlowMod** messages will be needed to fully configure a switch. However, OpenFlow is designed to give switches a lot of latitude to enable efficient processing, often at the expense of programmability and understandability:

- Pattern semantics. As discussed in preceding sections, the semantics of flow tables is non-trivial: patterns have implicit dependencies and flow tables can have multiple, overlapping entries. (The OpenFlow specification itself notes that scanning the table to find overlaps is expensive.) Therefore, it is up to the controller to avoid overlaps that introduce non-determinism.
- Packet reordering. Switches may reorder packets arbitrarily.
 For example, switches often have both a "fast path" that uses custom packet-processing hardware and a "slow path" that processes packets using a slower general-purpose CPU.
- No acknowledgments. Switches do not acknowledge when FlowMod messages are processed, except when errors occur. The controller can explicitly request acknowledgements by sending a barrier request after a FlowMod. When the switch has processed the FlowMod (and all other messages received before the barrier request), it responds with a barrier reply.
- Control message reordering. Switches may process control messages, including FlowMod messages, in any order. This is based on the architecture of switches, where the logical flow table is implemented by multiple physical tables working in parallel—each physical table typically only matches headers for one protocol. To process a rule with a pattern such as {nwSrc = 10.0.0.1, dlTyp = 0x800}, which matches headers across several protocols, several physical tables may need to be reconfigured, which takes longer to process than a simple pattern such as {dlDst = H2}.

Figure 7 defines the syntax and semantics of Featherweight Open-Flow, which faithfully models all of these behaviors. The rest of this section discusses the key elements of the model in detail.

6.2 Network Elements

Featherweight OpenFlow has four kinds of elements: switches, controllers, links between switches (carrying data packets), and links between switches and the controller (carrying OpenFlow messages). The semantics is specified using a small-step relation, with elements interacting by passing messages and updating their state non-deterministically.

Switches. A switch $\mathbb S$ comprises a unique identifier sw, a set of ports pts, and input and output packet buffers in_p and out_p . The buffers are multisets of packets tagged with ports, (pt, pk). In the input buffer, packets are tagged with the port on which they were received. In the output buffer, packets are tagged with the port on which they will be sent out. Since buffers are unordered, switches can process packets in any order. Switches also have a flow table, FT, which determines how the switch processes packets. As detailed in Section 4, the table is a collection of flow table entries, where each entry has a priority, pattern and, a multiset of output ports. Each switch also has a multiset of messages to

and from the controller, out_m and in_m . There are three kinds of messages from the controller:

- **PktOut** pt pk instructs the switch to emit packet pk on port pt.
- **FlowMod** δ instructs the switch to add or delete entries from its flow table. When δ is **Add** n pat act, a new entry is created, whereas **Del** pat deletes all entries that match pat exactly. In our model, we assume that flow tables on switches can be arbitrarily large. This is not the case for hardware switches, where the size of flow tables is often constrained by the amount of silicon used, and varies from switch-to-switch. It would be straightforward to modify our model to bound the size of the table on each switch.
- BarrierRequest n forces the switch to process all outstanding messages before replying with a BarrierReply n message.

Controllers. A controller $\mathbb C$ is defined by its local state σ , an input relation f_{in} , and an output relation f_{out} . The local state and these relations are application-specific, so Featherweight OpenFlow can be instantiated with any controller whose behavior can be modeled in this way. The f_{out} relation sends a message to a switch while f_{in} receives a message from a switch. Both relations update the state σ . There are two kinds of messages a switch can send to the controller:

- PktIn pt pk indicates that packet pk was received on pt and did not match any entry in the flow table.
- BarrierReply n indicates that sw has processed all messages up to and including a BarrierRequest n sent earlier.

Data links. A data link \mathbb{L} is a unidirectional queue of packets between two switch ports. To model bidirectional links we use symmetric unidirectional links. Featherweight OpenFlow does not model packet-loss in links and packet-buffers. It would be easy to extend our model so that packets are lost, for example, with some probability. Without packet loss, a packet traces paths from its source to its destinations (or loops forever). With packet loss, a packet traces a prefix of the complete path given by our model under ideal conditions.

Control links. A control link \mathbb{M} is a bidirectional link between the switch and the controller that contains a queue of controller messages for the switch and a queue of switch messages headed to the controller. Messages between the controller and the switch are sent and delivered in order, but may be processed in any order.

7. Verified Run-Time System

So far, we have developed a semantics for NetCore (Section 3), a compiler from NetCore to flow tables (Section 4), and a low-level semantics for OpenFlow (Section 6). To actually execute NetCore programs, we also need to develop a run-time system that installs rules on switches and prove it correct.

7.1 NetCore Run-Time System

There are many ways to build a controller that implements a Net-Core run-time system. A trivial solution is to simply process all packets on the controller. The controller receives input packets as **PktIn** messages, evaluates them using the NetCore semantics, and emits the outputs using **PktOut** messages.

Of course, we can do much better by using the NetCore compiler to actually generate flow tables and install those rules on switches using **FlowMod** messages. For example, given the following program,

 $dlDst = H1 \text{ and } not(dlTyp = 0x800) \Rightarrow \{1\}$

$$\begin{array}{lll} \text{Location} & loc ::= sw \times pt \\ \text{Located packet} & lp ::= loc \times pk \\ \text{Topology} & T & \in loc \rightharpoonup loc \end{array}$$

$$pg,T \vdash \{\!\lceil lp \rceil\!\rbrace \stackrel{lp}{\Longrightarrow} \{\!\lceil lp \rceil\!\rbrace$$

$$\begin{split} & \underbrace{lps' = \{ |(T(sw, pt_{out}), pk) \mid (pt_{out}, pk) \in \llbracket pg \rrbracket \ sw \ pt \ pk \} \}}_{pg, T \vdash \quad \{ |((sw, pt), pk) \} \ \uplus \ \{ | lp_1 \cdots lp_n \} \xrightarrow{(sw, pt, pk)}_{ps' \uplus \ \{ | lp_1 \cdots lp_n \} \}} \end{split}$$

Figure 8: Network semantics.

the compiler might generate the following flow table,

Priority	Pattern	Action
5	$\{\mathbf{dlDst} = \mathtt{H1}, \mathbf{dlTyp} = \mathtt{0x800}\}$	{ }
4	$\{dlDst = H1\}$	{ 1 }
3	*	{ }

and the controller would emit three FlowMod messages:

Add
$$5$$
 {dlDst = H1, dlTyp = 0x800} {||} Add 4 {dlDst = H1} {|1|} Add $3 \star \{||\}$

However, it would be unsafe to emit just these messages. As discussed in Section 6, switches can reorder messages to maximize throughput. This can lead to transient bugs by creating intermediate flow tables that are inconsistent with the intended policy. For example, if the $\mathbf{Add} \ 3 \ \star \{\|\}$ message is processed first, all packets will be dropped. Alternatively, if $\mathbf{Add} \ 4 \ \{\mathbf{dlDst} = \mathbf{H1}\} \ \{\|1\|\}$ is processed first, traffic that should be dropped will be incorrectly forwarded. Of the six possible permutations, only one has the property that all intermediate states either (i) process packets according to the program, or (ii) send packets to the controller (which can evaluate them using the program). Therefore, to ensure the switch processes the messages in order, the run-time system must intersperse **BarrierRequest** messages between **FlowMod** messages.

Network semantics. The semantics of NetCore presented in Section 3 defines how a program processes a single packet at a single switch at a time. But Featherweight OpenFlow models the behavior of an entire network of inter-connected switches with multiple packets in-flight. To reconcile the difference between these two, we need a *network semantics* that models the processing of all packets in the network. In this semantics (Fig. 8), the system state is a bag of in-flight located packets $\{lp\}$. At each step, the system:

- 1. Removes a located packet ((sw, pt), pk), from its state,
- Processes the packet according to the program to produce a new multiset of located packets,

$${lp_1 \cdots lp_n} = \llbracket pg \rrbracket \text{ sw pt } pk,$$

- 3. Transfers these packets to input ports, using the topology, $T(lp_1)\cdots T(lp_n),$ and
- 4. Adds the transferred packets to the system state.

Note that this approach to constructing a network semantics is not specific to NetCore: any hop-by-hop packet processing function could be used. Below, we refer to any semantics constructed in this way as a *network semantics*.

7.2 Run-Time System Correctness

Now we are ready to prove the correctness of the NetCore run-time system. However, rather than proving this directly, we instead develop a general framework for establishing controller correctness, and obtain the result for NetCore as a special case.

Bisimulation equivalence. The inputs to our framework are: (i) the high-level, hop-by-hop function the network is intended to implement, and (ii) the controller implementation, which is required to satisfy natural safety and liveness conditions. Given these parameters, we construct a weak bisimulation between the network semantics of the high-level function and an OpenFlow network instantiated with the controller implementation. This construction handles a number of low-level details once and for all, freeing developers to focus on essential controller correctness properties.

We prove a weak (rather than strong) bisimulation because Featherweight OpenFlow models the mechanics of packet processing in much greater detail than in the network semantics. For example, consider a NetCore program that forwards a packet pk from one switch to another, say SI to S2, in a single step. An equivalent Featherweight OpenFlow implementation would require at least three steps: (i) process pk at SI, move pk from the input buffer to the output buffer, (ii) move pk from SI's output buffer to the link to S2, and (iii) move pk from the link to S2's input buffer. If there were other packets on the link (which is likely!), additional steps would be needed. Moreover, pk could take an even more circuitous route if it is redirected to the controller.

The weak bisimulation states that the NetCore and Feather-Weight OpenFlow are indistinguishable modulo "internal" steps. Hence, any reasoning about the trajectory of a packet at the NetCore level will be preserved in Feather-Weight OpenFlow.

Observations. To define a weak bisimulation, we need a notion of observation (called an action in the π -calculus). We say that the NetCore network semantics observes a packet (sw, pt, pk) when it removes the packet from its state—i.e., just before evaluating it. Likewise, a Featherweight OpenFlow program observes a packet (sw, pt, pk) when it removes (pt, pk) from the input buffer on sw to process it using the FWD rule.

Bisimulation relation. Establishing a weak bisimulation requires exhibiting a relation \approx_{OF} between the concrete and abstract states with certain properties. We relate packets located in links and buffers in Featherweight OpenFlow to packets in the abstract network semantics. We elide the full definition of the relation, but describe some of its key characteristics:

- Packets (pt, pk) in input buffers inp on sw are related to packets ((sw, pt), pk) in the abstract state.
- Packets (pt, pk) in output buffers outp on sw are related to packets located at the other side of the link connected to pt.
- Likewise, packets on a data link (or contained in **PktOut** messages) are related to packets located at the other side of the data link (or the link connected to the port in the message).

Intuitively, packets in output buffers have already been processed and observed. The network semantics moves packets to new locations in one step whereas OpenFlow requires several more steps, but we must not be able to observe these intermediate steps. Therefore, after Featherweight OpenFlow observes a concrete packet pk (in the FWD rule), subsequent copies of pk must be related to packets at the ultimate destination.

The structure of the relation is largely straightforward and dictated by the nature of Featherweight OpenFlow. However, a few parts are application specific. In particular, packets at the controller and packets sent to the controller in **PktIn** messages may relate to the state in the network semantics in application-specific ways.

Abstract semantics. So far, we have focused on NetCore to build intuitions. But our bisimulation can be obtained for any controller

that implements a high-level packet-processing function. We now make this precise with a few additional definitions.

Definition 1 (Abstract Semantics). *An abstract semantics is defined by the following components:*

1. An abstract packet-processing function on located packets:

$$f(lp) = \{|lp_1 \cdots lp_n|\}$$

2. An abstraction function, $c:\sigma\to\{[lp]\}$, that identifies the packets the controller has received but not yet processed.

Note that the type of the NetCore semantics (Fig. 8) matches the type of the function above. In addition, because the NetCore controller simply holds the multiset of **PktIn** messages, the abstraction function is trivial. Given such an abstract semantics, we can lift it to a network semantics $\stackrel{lp}{\Longrightarrow}$ as we did for NetCore.

We say that an abstract semantics is *compatible* with a concrete controller implementation, consisting of a type of controller state σ , and input and output relations f_{in} and f_{out} , if the two satisfy the following conditions relating their behavior:

Definition 2 (Compatibility). *An abstract semantics and controller implementation are* compatible *if*:

- The controller ensures that all times packets are either (i) processed by switches in accordance with the packet-processing function or (ii) sent to the controller for processing;
- 2. Whenever the controller receives a packet,

$$(sw, PktIn \ pt \ pk, \sigma) \leadsto \sigma'$$

it applies the packet-processing function f to pk to get a multiset of located packets and adds them to its state

$$c(\sigma') = c(\sigma) \uplus f(pk)$$

3. Whenever the controller emits a packet,

$$\sigma \leadsto (sw, PktOut \ pt \ pk, \sigma')$$

it removes the packet from its state:

$$c(\sigma') = c(\sigma) \setminus \{|(sw, pt, pk)|\}$$

- 4. The controller eventually processes all packets (sw, pt, pk) in its state $c(\sigma)$ according to the packet-processing function, and
- 5. The controller eventually processes all OpenFlow messages.

The first property is essential. If it did not hold, switches could process packets contrary to the intended packet-processing relation. Proving it requires reasoning about the messages sent to the switches by the controller. In particular, because switches may reorder messages, barriers must be interspersed appropriately. The second and third properties relate the abstraction function c and the controller implementation. The fourth property requires the controller to correctly process every packet it receives. The fifth property is a liveness condition requiring the controller to eventually process every OpenFlow message. This holds in the absence of failures on the control link and the controller itself.

Given such a semantics, we show that our relation between abstract and Featherweight OpenFlow states and its inverse are weak simulations. This implies that the relation is a weak bisimulation, and thus that the two systems are weakly bisimilar.

Theorem 2 (Weak Bisimulation). For all compatible abstract semantics and controller implementations, all Featherweight Open-Flow states s and s', and all abstract states t and t':

• If $s \approx_{OF} t$ and $s \xrightarrow{(sw,pt,pk)} s'$, then there exists an abstract network state t'' such that $t \xrightarrow{(sw,pt,pk)} t''$ and $s' \approx_{OF} t''$, and

• If $s \approx_{OF} t$ and $t \xrightarrow{(sw,pt,pk)} t'$, then there exists a Featherweight OpenFlow state s'', and abstract network states s_i , s_i' such that

$$s \longrightarrow^* s_i \xrightarrow{(sw,pt,pk)} s'_i \longrightarrow^* s''$$

and
$$s'' \approx_{OF} t'$$
.

In this theorem, portions of the \approx_{OF} relation are defined in terms of the controller abstraction function, c supplied as a parameter. In addition, the proofs themselves rely on compatibility (Definition 2). Finally, we instantiate this theorem for the NetCore controller:

Corollary 1 (NetCore Run-Time Correctness). The network semantics of NetCore is weakly bisimilar to the concrete semantics of the NetCore controller in Featherweight OpenFlow.

8. Implementation and Evaluation

We have built a complete working implementation of the system described in this paper, including machine-checked proofs of each of the lemmas and theorems. Our implementation is available under an open-source license at the following URL:

Our system consists of 12 KLOC of Coq, which we extract to OCaml and link against two unverified components:

- A library to serialize OpenFlow data types to the OpenFlow wire format. This code is a lightly modified version of the Mirage OpenFlow library [16] (1.4K LOC).
- A module to translate between the full OpenFlow protocol and the fragment used in Featherweight OpenFlow (200 LOC).

We have deployed our NetCore controllers on real hardware and used them to build a number of useful network applications including host discovery, shortest-path routing, spanning tree, access control, and traffic monitoring. Using the union operator, it is easy to compose these modules with others to form larger applications.

NetCore at home. For the past month, we have used a NetCore controller to manage the home network of one of the authors. Home networks are small, but they are dynamic as devices regularly connect and disconnect. Our controller monitors the state of the network and dynamically reconfigures an OpenFlow-enabled wireless router to support the current devices. Because Featherweight OpenFlow does not model dynamic configurations, this controller's runtime system is necessarily unverified. However, it uses the NetCore compiler as a verified sub-system. In the future, we plan to enrich our OpenFlow model with additional features, including support for dynamic configurations.

Controller throughput. Controller throughput is important for the performance of SDNs. The CBench [26] tool quantifies controller throughput by flooding the controller with **PktIn** messages and measuring the time taken to receive **PktOut** messages in response. This is a somewhat crude metric, but it is still effective, since any controller must respond to **PktIn** messages. We used CBench to compare the throughput of our verified controller with our previous unverified NetCore controller, written in Haskell, and with the popular POX and NOX controllers, written in Python and C++ respectively. To ensure that the experiment tested throughput and not the application running on it, we had each controller execute a trivial program that floods all packets. We ran the experiment on a dual-core 3.3 GHz Intel i3 with 8GB RAM with Ubuntu 12.04 and obtained the results shown in Fig. 9 (a).

Our unverified NetCore controller is significantly faster than our verified controller. We attribute this to (i) a more mature backend that uses an optimized library from Nettle [27] to serialize

	26,022 16,997 9,437 6,150	Verified NetCore Verified NetCore 10 ⁵ 10 ⁴ 10 ⁴ 20 40 60 80	Univerified NetCore 10 ⁵ 10 ⁴ 10 ⁴ 10 ³ 20 40 60 80
		Time (seconds)	Time (seconds)
S1 S2 S3 S4 S5 S6		MicroFlow 10 ⁵ 10 ⁷ 10 ⁸ 20 40 60 80	PacketOut 10 ⁵ 10 ⁵ 10 ⁴ 20 40 60 80

Figure 9: Experiments: (a) controller throughput results; (b) control traffic topology; (c) control traffic results.

messages, and (ii) Haskell's superior multicore support, which the controller exploits heavily. However, despite being slower than the original NetCore, the new controller is still fast enough to be useful—indeed, it is faster than the popular POX controller (although POX is not tuned for performance). We plan to optimize our controller to improve its performance in the future.

(b)

Control traffic. Another key factor that affects SDN performance is the amount of traffic that the controller must handle. This metric measures the effectiveness of the controller at compiling, optimizing, and installing forwarding rules rather than processing packets itself. To properly assess a controller on these points, we need a more substantial application than "flood all packets." Using Net-Core, we built an application that computes shortest path forwarding rules as well as a spanning tree for broadcast. We ran this program on the six-switch Waxman topology shown in Fig. 9 (b), with two hosts connected to each switch.

In the experiment, every host sent 10 ICMP (ping) broadcast packets along the spanning tree, and received the replies from other hosts along shortest path routes. We used Mininet [9] to simulate the network and collected traffic traces using tcpdump. The total amount of network traffic during the experiment was 372 Kb.

We compared our Verified NetCore controller to several others: a (verified) "PacketOut" controller that never installs forwarding rules and processes all packets itself; our previous "Unverified NetCore" controller, written in Haskell; and a reactive "MicroFlow" controller [7] written in Haskell. The results of the experiment are shown in Fig. 9 (c). The graphs plot time-series data for every controller, showing the amount of control traffic in each one-second interval. Note that the y axis is on a logarithmic scale.

In the plot for our Verified NetCore controller, there is a large spike in control traffic at the start of the experiment, where the controller sends messages to install the forwarding rules generated from the program. Additional control traffic appears every 15 seconds; these messages implement a simple keep-alive protocol between the controller and switches. The Unverified NetCore controller uses the same compilation and run-time system algorithms as our verified controller, so its plot is nearly identical. The MicroFlow controller installs individual fine-grained rules in response to individual traffic flows rather than proactively compiling complete flow tables. Accordingly, its plot shows that there is much more control traffic than for the two NetCore controllers. The graph shows how traffic spikes when multiple hosts respond simultaneously to an ICMP broadcast. The fourth plot shows the behavior of

the PacketOut controller. Because this controller does not install any forwarding rules on the switches, all data traffic flows to the controller and then back into the network.

(c)

Although these results are preliminary, we believe they demonstrate that the performance of our verified NetCore controller can be competitive with other controllers. In particular, our verified controller generates the same flow tables and handles a similar amount of traffic as the earlier unverified NetCore controller, which was written in Haskell. Moreover, our system is not tuned for performance. As we optimize and extend our system, we expect that its performance will only improve.

9. Related Work

Verification technology has progressed dramatically in the past decades, making it feasible to prove useful theorems about real systems including databases [18], compilers [15], and even whole operating systems [14]. Compilers have been particularly fruitful targets for verification efforts [11]. Most prominently, the CompCert compiler translates programs in a large subset of C to PowerPC, ARM, and x86 executables [15]. The Verified Software Toolchain project provides machine-checked infrastructure for connecting properties obtained by program analysis to guarantees at the machine level [2]. Rocksalt verifies a tool for analyzing machine code against a detailed model of x86 [23]. Another system, Bedrock provides rich Coq libraries for verifying low-level programs [5]. Much earlier, a compiler for a Pascal-like language was formalized and verified as a part of the CLInc stack [31]. Significant portions of many other compilers have been formalized and verified, including the LLVM intermediate representation [33], the F* typechecker [25], and an extension of CompCert with garbage collection [20]. Our work is inspired by all of these efforts, but is the first to tackle the challenge of building a verified SDN controller.

Over the past few years, a number of researchers have proposed high-level programming languages for controlling networks, including COOLAID [4], FML [10], Frenetic [7], NetCore [22], and PANE [6]. This work uses NetCore [22] as a high-level network programming language. NetCore's original semantics was defined in terms of handwritten proofs and a complex abstract machine while we use machine-checked proofs and Featherweight Open-Flow. In proving our compiler and run-time system correct, we discovered several bugs in the unverified NetCore compiler and runtime. A portion of the PANE compiler was formalized in Coq, but since the proof did not model several subtleties of flow tables, the

compiler still had bugs. Unlike our system, PANE does not model or verify any portion of its run-time system. We used some of the PANE proofs during early development of our system. Lastly, Mirage [16], a language for writing cloud applications, includes an SDN interface. Our OpenFlow serializers are based on Mirage's.

Formally Verifiable Networking (FVN) [28] is a platform for synthesizing protocol implementations from formal specifications (though the synthesizer is unverified). Our work attacks the problem of generating and deploying correct network-wide configurations, rather than building distributed routing protocols. We use formal methods to build compilers, shifting the need for expertise with formal methods away from programmers.

Xie et al. introduced techniques for statically analyzing the reachability properties of networks [29]. A number of tools for verifying network configurations have been built using these techniques, including Header Space Analysis [12], Anteater [17], and VeriFlow[13]. These tools check whether the installed network rules have properties specified by the programmer. Our system guarantees that the generated network rules preserve the properties of the input program, enabling higher-level verification.

NICE [3] uses model-checking and symbolic execution to find bugs in OpenFlow controllers written in Python. Portions of our Featherweight OpenFlow model are inspired by the bugs discovered in NICE. Automatic Test Packet Generation [32] analyzes network configurations and constructs packets to achieve complete configuration testing coverage. Retrospective Causal Inference [24] detects minimal input sequences to induce bugs in SDN systems.

10. Conclusions

This paper presents a new foundation for network reasoning: a detailed model of OpenFlow, formalized in the Coq proof assistant, and a machine-verified compiler and run-time system for the Net-Core programming language. Our main result is a general framework for establishing controller correctness that reduces the proof obligation to a small number of safety and liveness properties. In the future, we plan to develop program logics for network programs, extend Featherweight OpenFlow with additional features not included in our current core calculus, and improve the engineering aspects of our system.

Acknowledgements. The authors wish to thank Andrew Ferguson and the PLDI reviewers for many helpful comments, and members of the Mirage project, especially Anil Madhavapeddy, for developing the OCaml CStruct and OpenFlow libraries. This work is supported in part by the NSF under grant CNS-1111698, the ONR under award N00014-12-1-0757, and by a Google Research Award.

References

- E. Al-Shaer and S. Al-Haj. FlowChecker: Configuration analysis and verification of federated OpenFlow infrastructures. In SafeConfig, 2010.
- [2] A. W. Appel. Verified software toolchain. In ESOP, 2011.
- [3] M. Canini, D. Venzano, P. Perešíni, D. Kostić, and J. Rexford. A NICE way to test OpenFlow applications. In NSDI, 2012.
- [4] X. Chen, Y. Mao, Z. M. Mao, and J. van der Merwe. Declarative configuration managaement for complex and dynamic networks. In *CoNEXT*, 2010.
- [5] A. Chlipala. Mostly-automated verification of low-level programs in computational separation logic. In *PLDI*, 2011.
- [6] A. D. Ferguson, A. Guha, C. Liang, R. Fonseca, and S. Krishnamurthi. Hierarchical policies for software defined networks. In *HotSDN*, 2012.
- [7] N. Foster, R. Harrison, M. J. Freedman, C. Monsanto, J. Rexford, A. Story, and D. Walker. Frenetic: A network programming language. In *ICFP*, 2011.

- [8] P. Gill, N. Jain, and N. Nagappan. Understanding network failures in data centers: measurement, analysis, and implications. In SIGCOMM, 2011.
- [9] N. Handigol, B. Heller, V. Jeyakumar, B. Lantz, and N. McKeown. Reproducible network experiments using container-based emulation. In *CoNEXT*, 2012.
- [10] T. L. Hinrichs, N. S. Gude, M. Casado, J. C. Mitchell, and S. Shenker. Practical declarative network management. In WREN, 2009.
- [11] T. Hoare. The verifying compiler: A grand challenge for computing research. JACM, 50(1):63–69, Jan 2003.
- [12] P. Kazemian, G. Varghese, and N. McKeown. Header space analysis: Static checking for networks. In NSDI, 2012.
- [13] A. Khurshid, X. Zou, W. Zhou, M. Caesar, and P. B. Godfrey. Verifiow: Verifying network-wide invariants in real time. In NSDI, 2013.
- [14] G. Klein, K. Elphinstone, G. Heiser, J. Andronick, D. Cock, P. Derrin, D. Elkaduwe, K. Engelhardt, R. Kolanski, M. Norrish, T. Sewell, H. Tuch, and S. Winwood. sel4: Formal verification of an OS kernel. In SOSP, 2009.
- [15] X. Leroy. Formal verification of a realistic compiler. CACM, 52(7): 107–115, Jul 2009.
- [16] A. Madhavapeddy, R. Mortier, C. Rotsos, D. Scott, B. Singh, T. Gazagnaire, S. Smith, S. Hand, and J. Crowcroft. Unikernels: Library operating systems for the cloud. In ASPLOS, 2013.
- [17] H. Mai, A. Khurshid, R. Agarwal, M. Caesar, P. B. Godfrey, and S. T. King. Debugging the data plane with Anteater. In SIGCOMM, 2011.
- [18] G. Malecha, G. Morrisett, A. Shinnar, and R. Wisnesky. Towards a verified relational database management system. In *POPL*, 2010.
- [19] A. Markopoulou, G. Iannaccone, S. Bhattacharyya, C.-N. Chuah, Y. Ganjali, and C. Diot. Characterization of failures in an operational IP backbone network. *IEEE/ACM Transactions on Networking*, 16(4): 749–762, Aug 2008.
- [20] A. McCreight, T. Chevalier, and A. Tolmach. A certified framework for compiling and executing garbage-collected languages. In *ICFP*, 2010
- [21] 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.
- [22] C. Monsanto, N. Foster, R. Harrison, and D. Walker. A compiler and run-time system for network programming languages. In POPL, 2012.
- [23] G. Morrisett, G. Tan, J. Tassarotti, J.-B. Tristan, and E. Gan. RockSalt: Better, faster, stronger SFI for the x86. In *PLDI*, 2012.
- [24] R. C. Scott, A. Wundsam, K. Zarifis, and S. Shenker. What, Where, and When: Software Fault Localization for SDN. Technical Report UCB/EECS-2012-178, EECS Department, University of California, Berkeley, 2012.
- [25] P.-Y. Strub, N. Swamy, C. Fournet, and J. Chen. Self-certification: Bootstrapping certified typecheckers in F* with Coq. In POPL, 2012.
- [26] A. Tootoonchian, S. Gorbunov, Y. Ganjali, M. Casado, and R. Sher-wood. On controller performance in software-defined networks. In HotICE, 2012.
- [27] A. Voellmy and P. Hudak. Nettle: Functional reactive programming of OpenFlow networks. In PADL, 2011.
- [28] A. Wang, L. Jia, C. Lio, B. T. Loo, O. Sokolsky, and P. Basu. Formally verifiable networking. In *HotNets*, 2009.
- [29] G. G. Xie, J. Zhan, D. A. Maltz, H. Zhang, A. G. Greenberg, G. Hjálmtýsson, and J. Rexford. On static reachability analysis of IP networks. In *INFOCOM*, 2005.
- [30] Z. Yin, M. Caesar, and Y. Zhou. Towards understanding bugs in open source router software. In SIGCOMM CCR, 2010.
- [31] W. Young. Verified compilation in micro-Gypsy. In TAV, 1989.
- [32] H. Zeng, P. Kazemian, G. Varghese, and N. McKeown. Automatic test packet generation. In *CoNEXT*, 2012.
- [33] J. Zhao, S. Nagarakatte, M. M. Martin, and S. Zdancewic. Formalizing the LLVM intermediate representation for verified program transformations. In *POPL*, 2012.