

RICE UNIVERSITY

**Controlling Race Conditions in OpenFlow to  
Accelerate Application Verification and Packet  
Forwarding**

by

**Xiaoye Sun**

A THESIS SUBMITTED  
IN PARTIAL FULFILLMENT OF THE  
REQUIREMENTS FOR THE DEGREE

**Master of Science**

APPROVED, THESIS COMMITTEE:



Dr. T. S. Eugene Ng, Advisor, Chair  
Associate Professor,  
Computer Science and Electrical and  
Computer Engineering



Dr. Edward W. Knightly  
Professor,  
Electrical and Computer Engineering



Dr. Lin Zhong  
Associate Professor,  
Electrical and Computer Engineering

Houston, Texas

October, 2014

## ABSTRACT

### Controlling Race Conditions in OpenFlow to Accelerate Application Verification and Packet Forwarding

by

Xiaoye Sun

OpenFlow is a Software Defined Networking (SDN) protocol that is being deployed in critical network systems. SDN application verification takes an important role in guaranteeing the correctness of the application. Through our investigation, we discover that application verification can be very inefficient under the OpenFlow protocol since there are many race conditions between the data packets and control plane messages. Furthermore, these race conditions also increase the control plane workload and packet forwarding delay. We propose *Attendre*, an OpenFlow extension, to mitigate the ill effects of the race conditions in OpenFlow networks. We have implemented *Attendre* in NICE (a model checking verifier), Open vSwitch (a software virtual switch) and NOX (an OpenFlow control platform). Experiments show that *Attendre* can reduce verification time by several orders of magnitude, and can significantly reduce TCP connection setup time.

# Contents

Abstract	ii
List of Illustrations	vi
List of Tables	vii
<b>1 Introduction</b>	<b>1</b>
1.1 Identify race conditions in an OpenFlow network . . . . .	3
1.2 Describe the ill effects of the race conditions . . . . .	3
1.3 Control and mitigate the ill effects of race conditions with Attendre .	4
1.4 Implement Attendre and quantify the benefits of Attendre . . . . .	4
<b>2 OpenFlow Basics</b>	<b>6</b>
<b>3 Races in OpenFlow and their Ill Effects</b>	<b>8</b>
3.1 Race Conditions in OpenFlow . . . . .	8
3.1.1 <b>Type 1</b> . . . . .	10
3.1.2 <b>Type 2</b> . . . . .	11
3.1.3 <b>Type 3</b> . . . . .	12
3.2 Ill Effects Due to the Race Conditions . . . . .	12
3.2.1 Increased complexity in SDN application verification . . . . .	12
3.2.2 Increased forwarding delay of the packet . . . . .	13
3.2.3 Increased processing workload on the switches and the controller	14
3.3 A Straw-man Solution to the Type 2 Race in OpenFlow . . . . .	14
<b>4 Attendre</b>	<b>15</b>

4.1	Overview of Attendre . . . . .	15
4.2	Attendre in simple examples . . . . .	16
4.2.1	<b>1 switch &amp; 2 packets</b> . . . . .	16
4.2.2	<b>2 switches &amp; 1 packet</b> . . . . .	17
4.3	Mechanisms in Attendre – how does Attendre solve the ill effect of race conditions in OpenFlow? . . . . .	18
4.3.1	Addressing type-1 race – Buffer the packets at the first switch	19
4.3.2	Addressing type-2 race – Buffer the packets at the following hop switches . . . . .	21
4.3.3	Addressing type-3 race – Gurantee one packet-in for each new flow . . . . .	23
4.3.4	Handling multiple flow-mods to a single switch – Message Set	23
4.3.5	Handling flow update – MS version . . . . .	24
4.3.6	Handling multiple flow tables . . . . .	24
<b>5</b>	<b>Implementation of Attendre</b>	<b>25</b>
5.1	Implementation in Open vSwitch . . . . .	25
5.2	Implementation in NOX . . . . .	26
5.3	Implementation in NICE . . . . .	26
<b>6</b>	<b>Evaluation</b>	<b>27</b>
6.1	Benefits to Model Checking Verification . . . . .	27
6.1.1	Route Flow . . . . .	27
6.1.2	MAC Learning Switch . . . . .	29
6.1.3	Energy-Efficient Traffic Engineering . . . . .	29
6.2	Benefits to TCP connection setup . . . . .	31
6.2.1	Experiment Setup . . . . .	31
6.2.2	Results . . . . .	34
6.3	Benefits to Forwarding Delay when Controller and Network Delays Vary	36

6.3.1	Delay Model in OpenFlow Network . . . . .	37
6.3.2	Delay Formulae . . . . .	38
6.3.3	Experiment Setup . . . . .	40
6.3.4	Numerical Analysis Results . . . . .	42
<b>7</b>	<b>Related Work</b>	<b>43</b>
7.1	Previous work related to SDN verification . . . . .	43
7.2	Previous work aiming at reducing the first packet forwarding delay and control plane workload in SDN . . . . .	45
<b>8</b>	<b>Conclusion</b>	<b>47</b>
	<b>Bibliography</b>	<b>48</b>

# Illustrations

3.1	Example scenarios showing race conditions in the OpenFlow protocol	9
3.2	Partial state diagrams for the scenarios in Fig. 3.1 . . . . .	10
4.1	Partial state diagrams for Attendre . . . . .	16
4.2	The propagation of Attendre command . . . . .	20
6.1	NICE verification results under various scenarios for OpenFlow and Attendre . . . . .	28
6.2	CDF of TCP connection establishment time with different hops and flow inter-arrival time . . . . .	32
6.3	Delay Analysis . . . . .	41

# Tables

3.1	Notations used in this thesis . . . . .	9
6.1	Number of State Transitions and Execution Time in MAC Learning Application . . . . .	30
6.2	Number of State Transitions and Execution Time in Energy-Efficient Traffic Engineering Application . . . . .	31
6.3	Delay Formulae Notation . . . . .	37

# Chapter 1

## Introduction

The OpenFlow protocol enables the programmatic and centralized control of a network's behavior. Software application programs running on a centralized OpenFlow controller can observe virtually any aspects of an OpenFlow-enabled network and correspondingly dictate the behavior of the network by sending command messages to directly configure packet processing rules on the OpenFlow-enabled network switches. This centralized, user-programmable control plane architecture is generally referred to as Software-Defined Networking (SDN). OpenFlow is currently the *de facto* standard protocol for SDN in the industry [1].

One of the most important usages of SDN is in fine-grained network security policy enforcement in enterprise networks to combat the unauthorized resource access threat [2–7]. To realize fine-grained network security policies, a critical feature supported by OpenFlow is the so-called packet-in messages. Using this feature, a traffic flow that is not known to the network switches is directed, via packet-in messages, to the centralized controller for processing. There, an application program can reactively make sophisticated decisions about how this traffic flow should be handled (e.g. authorized and routed, denied and filtered, redirected, monitored, logged etc.) based on a user customizable policy.

In such usages of SDN, the correctness of the SDN controller application programs is paramount. The SDN controller applications are responsible for both realizing the security policy and for recovering the network from failures [8] in switches, load



balancers, etc. [9]. It is therefore critical to rigorously verify the correctness of SDN controller applications.

Among all automated SDN verification approaches proposed (a broader discussion of the various approaches can be found in Chapter 7), model checking is the most powerful one, since it is capable of detecting the widest range of misbehaviors (e.g. forwarding loops, reachability failures, policy violations, lack of liveness) in SDN applications even before these applications are deployed. In model checking, the entire network, including the switches, the hosts and the controller running SDN applications, is modeled as a state machine. The verifier explores the entire state space and checks for network property invariant violations by simulating every possible event ordering inside the entire network.

Unfortunately, controlling the network in a reactive way using the packet-in message makes the model checking based application verification highly inefficient. The underlying reason is that the data packets and the OpenFlow messages exhibit many race conditions, and each race condition exacerbates the number of state transitions that the model checker needs to explore. We further discover that the race condition also increases packet forwarding delay and the processing overhead on switches and the controller. These problems greatly limit the practicality of SDN application verification as a means to ensure network correctness, and lower the performance of an SDN network unnecessarily.

This thesis is the first study to systematically analyze the origins and impacts of such race conditions in the OpenFlow protocol and to explore new protocol designs to mitigate the problems they cause. We make the following contributions:\*

---

\*The following discussion assumes that the reader has a basic understanding of the OpenFlow protocol and its terminologies. See Chapter 2 for a brief tutorial.

## 1.1 Identify race conditions in an OpenFlow network

We present and illustrate the following three types of race conditions in detail with two different network scenarios.<sup>†</sup> The first two types of races exist between a data packet and a command message. The third race exists between two command messages.

- At the switch sending out a packet-in message to the controller, the following packets in the same flow of the packet in the packet-in message race against the command messages sent to the switch for processing this flow.
- When the controller receives a packet-in and decides to send flow-mod messages to multiple switches for the data flow, at the switch not sending the packet-in but receiving the flow-mods, the data packets in the flow race against the flow-mods.
- Command messages sent to the same switch for the same flow race against each other.

## 1.2 Describe the ill effects of the race conditions

Using state transition diagrams, we point out the ill effects of the races in model checking based OpenFlow application verification. We also describe other ill effects, like increased packet forwarding delay and increased workload on the switches and the controller.

---

<sup>†</sup>Note that other types of race conditions exist in OpenFlow. However, they are less common and we do not present them due to space limitation.

### **1.3 Control and mitigate the ill effects of race conditions with Attendre**

We present Attendre, an extension of the OpenFlow protocol to control the race conditions so as to mitigate the ill effects of races. By controlling these races, we mean that the messaging behavior of a switch is as if the outcomes of the races are deterministic. By controlling the outcomes of the race conditions, Attendre eliminates many unnecessary network states as well as many unnecessary control plane messages, which are the root causes for inefficient verification and increased forwarding delay. We present the changes to the OpenFlow protocol and the changes to OpenFlow switches in detail.

### **1.4 Implement Attendre and quantify the benefits of Attendre**

We have implemented Attendre in NICE [10] (a model checking based OpenFlow application verifier), in Open vSwitch [11] (a software OpenFlow switch) and in NOX [12] (an OpenFlow controller platform). With these implementations, we demonstrate that comparing with original OpenFlow, Attendre could reduce the verification execution time by several orders of magnitude; we also show that Attendre could greatly reduce the TCP connection setup time. To theoretically understand the forwarding delay in an OpenFlow network with and without Attendre, we derive a forwarding delay model and numerically analyze the benefits on forwarding delay that Attendre can achieve under various controller processing delays and controller-to-switch link delays.

The rest of this thesis is organized as follows. In Chapter 2, we briefly review

the concepts and terminologies in OpenFlow. In Chapter 3, we illustrate the race conditions in OpenFlow and their ill effects. We present the Attendre mechanisms in Chapter 4. Chapter 5 presents the implementation of Attendre. Chapter 6.1, 6.2 and 6.3 show the experimental results that quantify the benefits of Attendre on verification and packet forwarding delay. We present related work in Chapter 7. We conclude and point to future directions in Chapter 8.

## Chapter 2

# OpenFlow Basics

SDN separates the data plane and the control plane of a network. The data plane, usually the switch ASIC chip in a hardware switch platform or the kernel module in a software virtual switch, is responsible for fast packet processing, such as forwarding or dropping packets and modifying packet headers. The control plane, on the other hand, configures the packet processing rules on data plane and handles the packets that the data plane don't know how to process, e.g., mismatching packets.

OpenFlow is a *de facto* standard SDN protocol the control plane uses to configure the network data plane. More specifically, configuring the data plane is achieved by sending OpenFlow command messages between a logically centralized controller and the switches. An OpenFlow switch has multiple pipelined ***flow tables***, each of which contains multiple flow table ***entries***.

The entry includes ***match fields***, a ***priority value*** and ***actions***. The match fields are defined in the OpenFlow specification. They consist of particular fields of the packet header, like MAC address, IP address and TCP port, and the ingress port of the packet. A packet will start matching at the flow table 0 and collect ***actions*** in the matched entries to the ***action set*** throughout the matching process. The action set is an ordered list associated to each packet and it contains the actions that will be applied to the packet at the end of the matching.

An OpenFlow switch communicates with the controller over a secured TCP channel. In OpenFlow, if a packet does not match with a flow table, the packet will

buffer\* at the switch and will be associated to a *buffer ID*, which can be used to retrieve the packet from the buffer. Then, a *packet-in* message containing the *match fields* of the packet, buffer ID and the *table ID* of the flow table at which the mismatch happens will be sent from the switch to the controller. By processing the packet-in, the controller could add one or more entries to the specified flow tables of switches by sending *flow-mod* messages to them. The controller could also issue a *packet-out* message containing actions and the buffer ID in the packet-in message back to the switch. The switch will process this packet according to the actions in the packet-out message. The flow-mod could also carry a buffer ID, so that a packet-out message is combined with the flow-mod. We denote this flow-mod message as a *flow-mod/packet-out* message. Throughout the thesis, we do not differentiate flow-mod from flow-mod/packet-out in some cases, but the context should make it clear.

---

\*The switch can also send the entire packet to the controller together with the packet-in message.

## Chapter 3

### Races in OpenFlow and their Ill Effects

In this Chapter, we reveal different types of races and their implications in detail. The following discussions assume that before the initial packet comes to the first switch, the flow table entries that match the packet are absent on the switches. This is a common case in a network that implements a fine-grained flow authentication and security management. For example, Ethane [13] and Resonance [7] only install rules for a flow when the host starts a new connection or the host sending the packet just join the network, since each new flow or host address needs to be authenticated by the controller before it traverses the network. The notations used are shown in Table 3.1.

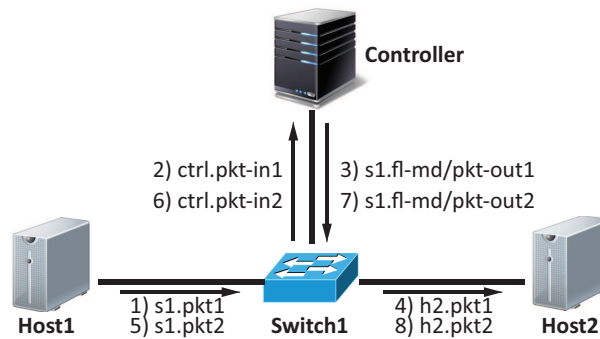
The controller application we use here is “route flow” unless otherwise mentioned. This application installs forwarding entries to all the switches on the forwarding path of each packet. In detail, the application issues a flow-mod/packet-out message to the switch that sent the packet-in message (first switch) and flow-mod messages to the following switches that are on the rest of the forwarding path.

#### 3.1 Race Conditions in OpenFlow

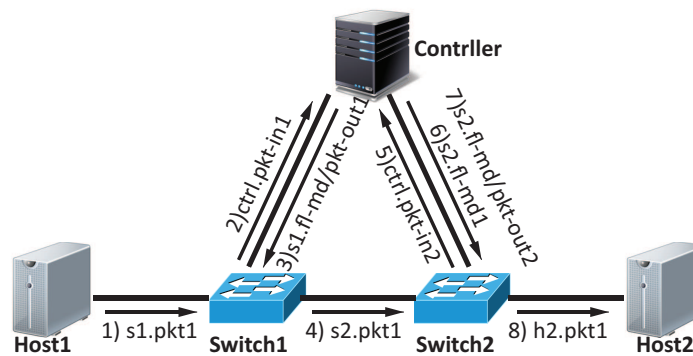
We have identified three races prevalent in an OpenFlow network. The network topologies and the state diagrams in Fig. 3.1 and 3.2 are shown here to help explain the races. Only part of the complete state diagrams are shown for brevity.

Table 3.1 : Notations used in this thesis

Notation	Meaning
$hn$	$n$ th host
$sn$	$n$ th switch
ctrl	OpenFlow controller
$pktn$	$n$ th data packet
$pkt-inn$	$n$ th packet-in message
$fl-mdn$	$n$ th flow-mod message
$fl-md/pkt-outn$	$n$ th flow-mod/packet-out message



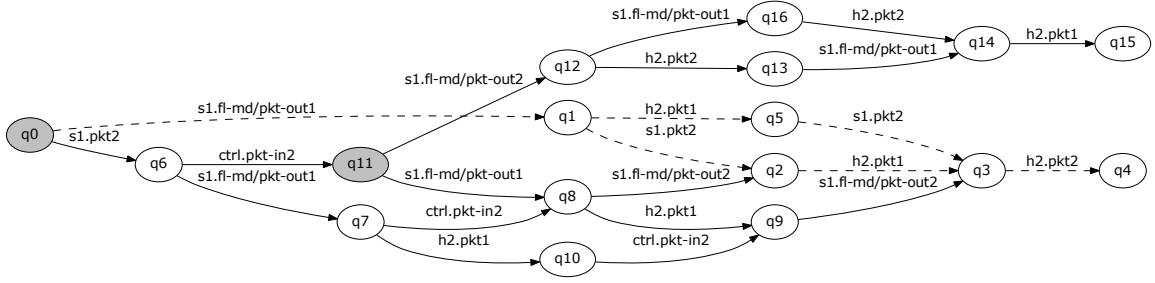
(a) Network diagram for one switch and two packets



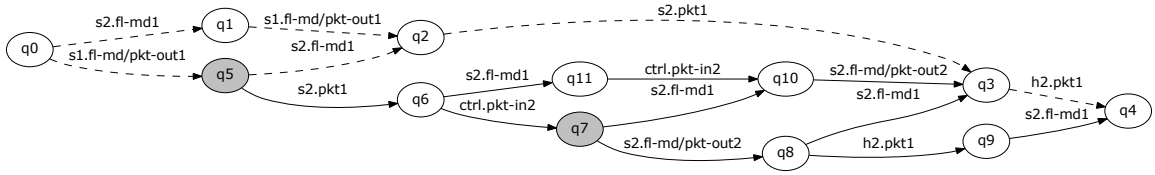
(b) Network diagram for two switches and one packet

Figure 3.1 : Example scenarios showing race conditions in the OpenFlow protocol





(a) OpenFlow state diagram depicting race type 1 and 3



(b) OpenFlow state diagram depicting race type 2 and 3

Figure 3.2 : Partial state diagrams for the scenarios in Fig. 3.1

### 3.1.1 Type 1

This race happens at the first switch and the race is between the processing of the non-initial data packets and the processing of the flow-mod messages containing the flow table entries addressing those packets. Fig. 3.1(a) is used to illustrate this type of race with one switch in the network and two packets sent from  $h1$  to  $h2$ .  $pkt1$  mismatches the flow table on  $s1$  (Label 1), so  $s1$  sends  $pkt-in1$  to  $ctrl$ . Label 2 here means  $ctrl$  has processed  $pkt-in1$  for  $pkt1$ , which is represented by state  $q_0$  in Fig. 3.2(a). Having processed  $pkt-in1$ ,  $ctrl$  sends  $fl-md/pkt-out1$  to  $s1$ . Meanwhile,  $h1$  also sends  $pkt2$  to  $h2$ , racing against  $fl-md/pkt-out1$ .

There are two possible outcomes depending on whether  $s1$  processes  $fl-md/pkt-out1$  (Label 3) or  $pkt2$  (Label 5) first. If Label 3 happens before Label 5,  $s1$  will have a flow table entry for  $pkt2$ . When  $pkt2$  arrives  $s1$ , it will be processed accord-

ing to the just installed entry, so that the state diagram takes the transitions from  $q_0 \rightarrow q_1 \rightarrow \dots \rightarrow q_4$ . However, consider the other case where  $pkt2$  is processed by  $s1$  before  $fl-md/pkt-out1$ . In this case,  $ctrl$  will process a redundant packet-in message from  $s1$  for  $pkt2$ , which means Label 6 will take place according to Fig. 3.1(a). The state diagram in Fig. 3.2(a) takes the transitions from  $q_0 \rightarrow q_6 \rightarrow \dots \rightarrow q_4/q_{15}$ . In addition,  $pkt2$  will be processed by  $s1$  only after the switch processes  $fl-md/pkt-out2$  (Label 7). Thus,  $pkt2$  suffers an extra delay before being forwarded.

### 3.1.2 Type 2

This race is experienced between data packets and flow-mod messages addressing these packets but not at the switch sending the packet-in message. Fig. 3.1(b) is used to illustrate this type of race with two switches and one injected packet. Fig. 3.2(b) is part of the state diagram for this case. Label 2 and state  $q_0$  in Fig. 3.1(b) and 3.2(b) respectively depict the same event as described in the type 1 race. Since there are two switches in the network,  $ctrl$  will also send  $fl-md1$  to  $s2$  to install a forwarding path for  $pkt1$ . After  $s1$  processes  $fl-md/pkt-out1$ ,  $pkt1$  will be sent to  $s2$  (Label 4). Meanwhile, if  $s1$  processes  $fl-md/pkt-out1$  (Label 3) before  $s2$  processes  $fl-md1$  (Label 6), represented by  $q_5$  in Fig. 3.2(b),  $fl-md1$  is still on its way to  $s2$ , racing against  $pkt1$ .

If Label 6 comes before Label 4,  $pkt1$  will be addressed by the just inserted flow table entry on  $s2$  by the state transitions from  $q_5 \rightarrow q_2 \rightarrow q_3 \rightarrow q_4$  in Fig. 3.2(b). Consider the case where Label 4 happens before Label 6.  $pkt1$  cannot match the flow table and a  $pkt-in2$  message will be sent by  $s2$  (Label 5). This  $pkt-in2$  is redundant since the forwarding rule has already been sent or is going to be sent to  $s2$  via Label 6 and it also introduces an extra delay to the forwarding of  $pkt1$ . Fig. 3.2(b) shows the

corresponding state transition for this case from  $q_5 \rightarrow q_6 \rightarrow \dots \rightarrow q_4$ .

### 3.1.3 Type 3

This race is experienced between different command messages issued by a multi-threaded controller for the same switch. This happens due to the controller threads competing for the same set of resources like processors or memories. In Fig. 3.1(b), this race happens at  $q_7$ , which represents the state where  $s2$  processes  $pkt1$  (Label 4) before it processes  $fl-md1$  (Label 6), so that  $s2$  sends  $pkt-in2$  and it has been processed by  $ctrl$  (Label 5). In this case, it is possible that the thread responsible for sending  $fl-md1$  is delayed by the controller. Thus, it is possible that  $s2$  will process  $fl-md/pkt-out2$  before  $fl-md1$  and vice versa. In the former case, the state diagram in Fig. 3.2(b) will take the state transitions  $q_7 \rightarrow q_8 \rightarrow \dots \rightarrow q_4$ . But, if  $s2$  processes  $fl-md1$  first, the state transitions  $q_7 \rightarrow q_{10} \rightarrow \dots \rightarrow q_4$  will be taken. The same type of race happens at  $q_{11}$  in Fig. 3.2(a). As we can see, this race also results in many more state transitions.

## 3.2 Ill Effects Due to the Race Conditions

### 3.2.1 Increased complexity in SDN application verification

These races results in many possible event orderings. In addition, the unnecessary command messages also add more possible events into the system. In the model checking based application verification, the verifier needs to explore every possible event ordering in the system\*. The state diagrams shown in Fig. 3.2 are for the

---

\*The verification we will focus on in this thesis is the model checking based OpenFlow application verification, which will be called as “verification” in short unless otherwise mentioned.

simplest scenarios, where there are no more than two switches or packets. However, when the network scales up and has more switches and more packets going through the network, these races could result in a state explosion problem. The hollow markers in Fig. 6.1 show the numbers of state transitions and the time during verification on the “route flow” application with various numbers of switches and packets. We use NICE [10], a model checking based OpenFlow application verifier to verify the system. We can see that when the number of switches or packets increases, the numbers of state transitions and the verification time increase at an alarming rate. It is important to note that NICE assumes a single-threaded controller, therefore the numbers are highly conservative. If a multi-threaded controller model is used in NICE, the type 3 race will also be introduced and will cause further explosion in the state space.

### 3.2.2 Increased forwarding delay of the packet

In OpenFlow, a packet that incurs a packet-in will be processed by that switch only after the switch has processed the corresponding packet-out message. The latency from the switch sending the packet-in to the switch receiving the packet-out is composed of the round trip time between the switch and the controller and the controller processing delay. In the worst case, this could happen at every hop on the path. However, the packet can be addressed by the switch as soon as the related command messages have been processed by the switch. Thus, an unnecessary packet-in that results from the race increases the forwarding delay of the packet significantly since in some cases the corresponding command messages would have been processed by the switch just after a packet-in has been sent to the controller.

### 3.2.3 Increased processing workload on the switches and the controller

The unnecessary packet-in, flow-mod and packet-out messages caused by the races increase the workload both on the switches and the controller. These redundant messages will further limit the scalability of an OpenFlow controller and also increase the switch processing time.

## 3.3 A Straw-man Solution to the Type 2 Race in OpenFlow

In the OpenFlow protocol, the application could use *barrier messages* to avoid the type 2 race. More specifically, the controller application could hold the flow-mod/packet-out for the first switch and send a *barrier request* message following each flow-mod for the subsequent switches. A switch receiving a barrier request will send a *barrier reply* back to the controller after it finishes processing all the command messages received before the barrier request. The controller application will send the flow-mod/packet-out back to the first switch after it receives all the barrier replies from the following switches. As a consequence, when the packets reach the following switches, the flow table entries for them have already been installed. Thus, the type 2 race is eliminated.

However, using barrier messages does not eliminate the type 1 and type 3 races. In addition, the waiting at the controller introduces delays to packet forwarding. We will evaluate the drawback of using barrier messages in Chapter 6.1, 6.2 and 6.3. In the following chapters, **OpenFlow-B** will be used to denote the scenarios where the OpenFlow application uses barrier message in the way we discussed above.

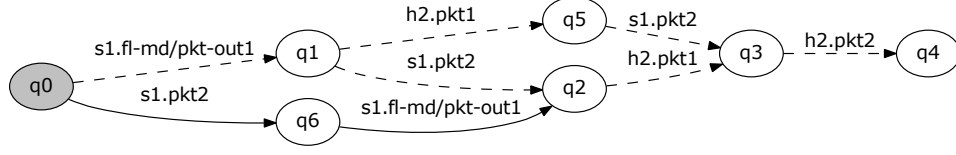
## Chapter 4

### Attendre

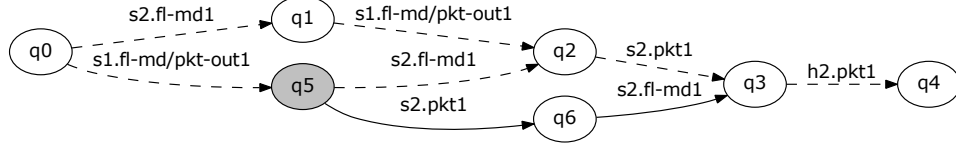
#### 4.1 Overview of Attendre

As we discussed in the Chapter 3, the race conditions in OpenFlow result in many ill effects. To mitigate these ill effects, we develop Attendre, which extends the OpenFlow protocol with a series of mechanisms. The key idea behind Attendre is that it controls the outcomes of the race conditions and effectively always makes the command message win the first two types of races.

In more detail, the switch avoids sending unnecessary packet-in messages, e.g., Label 6 in Fig. 3.1(a) and Label 5 in Fig. 3.1(b), by buffering the data packets at the switch until the switch processes the flow-mod and packet-out messages for the packets. (It may at first seem counter-intuitive, but the discussion in this chapter and the experiments in Chapter 6.3 and 6.2 will show that the packet buffering *decreases* the forwarding delay of a packet.) By doing so, each new flow only generates one packet-in message. It also avoids other unnecessary command messages, e.g., Label 7 in Fig. 3.1(a) and Label 7 in Fig. 3.1(b). To avoid buffering all the incoming packets causing the network to be stuck, the switch passes an Attendre Command (AC) to the next hop switch on the forwarding path of the packets. The AC tells the switch which flow it should buffer. After the switch processed the command messages, it releases the buffered packets.



(a) Attendre state diagram for Fig. 3.1(a)



(b) Attendre state diagram for Fig. 3.1(b)

Figure 4.1 : Partial state diagrams for Attendre

## 4.2 Attendre in simple examples

In this part, we use the two scenarios in Fig. 3.1 to explain how Attendre works and how Attendre reduce the cost in SDN application verification, packet forwarding delay and control plane overhead.

When a packet fails to match the flow table, if the switch is expecting to receive command messages addressing the packet, the switch should buffer the packet in its buffer. Fig. 4.1 gives the intuition on how Attendre works by showing the state diagrams for the network scenarios in Fig. 3.1. Dashed edges are the common state transitions shared between Fig. 3.2 and 4.1. Comparing the state diagrams for OpenFlow with and without Attendre, we can directly see that many state transitions are eliminated by Attendre.

### 4.2.1 1 switch & 2 packets

With Attendre, the state diagram of this case is changed from Fig. 3.2(a) to Fig. 4.1(a). The shaded  $q_0$  has the same meaning as that in Fig. 3.2(a). At  $q_0$ , where type 1 race

happens, if we take the branch where  $s1$  processes  $fl-md/pkt-out1$  before it processes  $pkt2$ ,  $pkt2$  will find an entry in the flow table, then  $pkt2$  will be addressed by the matched entry. The dashed edges in Fig. 4.1(a),  $q_0 \rightarrow q_1 \rightarrow \dots \rightarrow q_4$ , shows the state transitions of the above process. With Attendre, when  $pkt1$  is sent out via  $pkt-in1$ , the switch will extract the match fields of  $pkt1$  and know that the controller is processing the packets having the match fields of  $pkt1$ . The match fields are defined in the OpenFlow specification. They consist of particular fields of the packet header, like MAC address, IP address and TCP port, and the ingress port of the packet. When  $pkt2$  matching the match fields of  $pkt1$  arrives at  $s1$ , it will be buffered, which is represented by the state transition of  $q_0 \rightarrow q_6$ . Thus,  $s1$  and the controller no longer need to send or process  $pkt-in2$  and  $fl-md/pkt-out2$ , which are possible outcomes of the type 1 race. After  $s1$  processes  $fl-md/pkt-out1$ , which means that the entry for the packets has been installed in the flow table,  $pkt1$  and the buffered  $pkt2$  will be sent to  $h2$  by re-matching the flow table and take the actions in the matched entry. This process is represented by the state transitions from  $q_6 \rightarrow \dots \rightarrow q_4$ . Verification benefits from this buffering as the races and transitions related to the second packet-in will not appear in the state diagram. A switch will address the buffered packets as soon as the corresponding commands are obtained by the switch, this reduces the delay in processing the following packets.

#### 4.2.2 2 switches & 1 packet

The state diagram in Fig. 4.1(b) shows the state transitions of the scenario in Fig. 3.1(b) in Attendre. With Attendre, after the controller processes  $pkt-in1$ , the controller will inform  $s1$  that it will send  $fl-md1$  to  $s2$ . This information will be sent to  $s1$  together with  $fl-md/pkt-out1$ . The same information will then be piggybacked by  $pkt1$ , which



---

**Algorithm 1** Procedure for processing a data packet
 

---

```

1: START: SWITCH receives a data packet DP
2: if DP has encapsulated AC then
3:   # process the Attendre Command
4:   if SWITCH has not processed all the messages in the MS at AC.MSVersion then
5:     if there is an ENTRY has the same match fields, priority and table id as the AC then
6:       if AC.MSVersion > ENTRY.MSVersion then
7:         replace ENTRY.action with buffer action
8:       else if AC.MSVersion = ENTRY.MSVersion then
9:         append buffer action to ENTRY.action
10:      end if
11:     else
12:       add a buffer entry with AC
13:     end if
14:   end if
15: end if
16: # process the original data packet
17: while initial matching or goto next flow table action do
18:   if DP matches the current flow table then
19:     execute actions in the matched entry, i.e., write actions to the action set of DP
20:   else
21:     # process mismatched data packet
22:     add an entry with the match fields of DP, highest priority, buffer action to the current
23:     flow table
24:     send packet-in for DP to the controller
25:   return
26: end if
27: end while
28: execute the action set of DP
29: return

```

---

matches the entry inserted to the flow table by *fl-md/pkt-out1*. *s2* on receiving *pkt1* will know that the controller has sent or is about to send *fl-md1* addressing *pkt1*. *s2* will buffer *pkt1* in its buffer until *s2* has processed the expected *fl-md1*. After that, the buffered *pkt1* will match the flow table again and thereby be addressed by that entry. The corresponding state transition is  $q_5 \rightarrow q_6 \rightarrow q_3$ .

### 4.3 Mechanisms in Attendre – how does Attendre solve the ill effect of race conditions in OpenFlow?

Attendre provides a series of mechanisms to mitigate the ill effects of the race conditions in OpenFlow. These mechanisms involves minor changes to the processing

---

**Algorithm 2** Procedure for processing a flow-mod message
 

---

```

1: START: SWITCH receives a flow-mod message FM
2: if FM contains a packet, i.e., flow-mod/packet-out then
3:   if FM has an AC then
4:     append encap action to the action set of the packet in FM
5:   end if
6:   process the packet in FM
7: else
8:   if FM has an AC then
9:     append the encap action to the action fields of the entry in FM
10:  end if
11: end if
12: add/modify the entry in FM to the flow table
13: if FM is the last message in the MS at FM.MSVersion then
14:   remove the buffer action
15:   let the packets buffered by the entry to re-match the flow table
16: end if
17: return

```

---



---

**Algorithm 3** Procedure for processing a packet-out message
 

---

```

1: START: SWITCH receives a packet-out message PO
2: if PO has an AC then
3:   append encap action to the action set of the packet in PO
4: end if
5: process the packet in PO
6: if PO is the last message in the MS at PO.MSVersion then
7:   remove the buffer entry
8:   let the packets buffered by the entry to re-match the flow table
9: end if
10: return

```

---

procedure of the data packet, flow-mod and packet-out messages. The new processing algorithms are shown in Alg. 1, 2 and 3 respectively.

#### 4.3.1 Addressing type-1 race – Buffer the packets at the first switch

If the first packet of a flow fails to match the flow table, the switch sends a packet-in message to the OpenFlow controller. In the rest of the section, we will call the switch that sends the packet-in message as **packet-in switch** of the flow. In Attendre, the packet-in switch avoids sending redundant packet-in messages for the following packets in the same flow, e.g., Label 6 in Fig. 3.1(a), by buffering these packets.

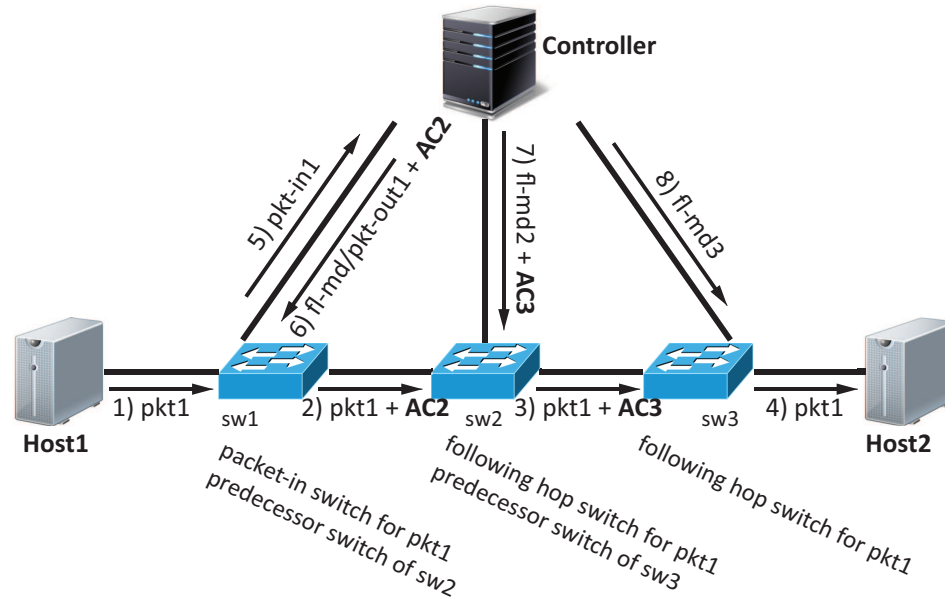


Figure 4.2 : The propagation of Attendre command

Attendre provides two mechanisms to support this function:

**The “*buffer*” action** – In OpenFlow, a packet is processed by a list of actions of the matched flow table entry. To make the packet buffering mechanism consistent with the OpenFlow protocol, we extend the protocol with a *buffer* action \*. This action stores the matched packet together with its action-set to the switch, assigns an unique buffer ID for each packet, and inserts the buffer ID to a queue associated with that entry. With the buffer ID, the switch could retrieve the packets after it processes all of the flow-mod and packet-out messages for that flow.

**Insert an entry with the *buffer* action** – To buffer the following packets of a flow, once the first packet mismatches the flow table at the packet-in switch, the switch inserts an entry to the mismatched flow table as shown by Line 21-24 in

---

\*The *buffer* action is an **apply-action**, which processes the packet immediately when the packet matches with the entry.

Alg. 1. This entry contains the match fields of the first packet, is assigned the highest priority value and is given the *buffer* action. By doing so, each new flow only creates one packet-in message at the packet-in switch; when the following packets arrive at the packet-in switch, the switch buffers them.

#### 4.3.2 Addressing type-2 race – Buffer the packets at the following hop switches

By saying **following hop switches** of a flow, we mean the switches on the forwarding path of the flow not including the packet-in switch. In Attendre, if a following hop switch has not processed the flow-mod messages for the flow, that switch also buffers the packets in that flow including the first packet. To achieve this function, the following hop switch needs to know which packets it should buffer before it starts to process the packets and buffers the packets until the switch processes the flow-mod messages for that flow. Attendre provides the following mechanisms to support this function:

**The Attendre command** – In Attendre, when the controller sends flow-mods to the following hop switches for a new flow, it uses the Attendre Command (AC) to tell each of the following hop switch which flow it should buffer. The controller sends the AC via a special mechanism which we discuss in the following paragraph. An AC includes match fields, a priority value, and a table ID. With these fields, when a following hop switch receives an AC, if it has not processed the flow-mod messages, the switch inserts an entry containing these fields and the *buffer* action. It is easy for the controller to generate the AC for each of the following hop switches since it is able to get the entries that will be installed on the switches by checking the flow-mod messages. For example, in the scenario in Fig. 3.1(b), the fields in the AC for *sw2*

are the same as the match fields, priority value and table ID of *fl-md1*.

**The *encap* action** – How the AC propagates in Attendre network is shown in Fig. 4.2. To avoid sending packet-ins to the controller during a flow update, a following hop switch should process the AC before the data packets in that flow, so that when the packets arrive at the switch, there is always an entry the packet can match with. This matched entry can be either a regular entry installed by the flow-mod or the *buffer* entry installed by the AC. In order to process the AC always before the packets in the flow, Attendre puts the AC into the data plane. See Label 2 and 3 in Fig. 4.2. More specifically, the AC for a switch is always encapsulated in the first data packet of the flow. When a following hop switch receives the AC from the data plane, it first takes out the AC, installs the *buffer* entry to the switch if its has not processed the flow-mod for the flow, and then processes the original data packet. (Line 3-15 in Alg. 1).

Again, to be consistent with the OpenFlow protocol, we add an *encap* action that encapsulates the AC to the data packet. This action will be removed immediately from the entry once it is used, so that the AC is only encapsulated to the first packet.  
†

In Attendre, the controller sends the AC for a following hop switch to its **predecessor switch** on the forwarding path. In this case, the AC is received by a switch via the control plane. The procedure of processing the AC received from the control plane is difference from processing the AC received from the data plane. If the predecessor switch is not the packet-in switch (*sw2* in Fig. 4.2), the AC is carried by the flow-mod that modifies the earliest flow table of the predecessor switch(Label 7 in

---

†The *encap* action is a Write-Action, which writes the *encap* action into the actions set of the packet. The switch executes the *encap* action right before the *output* action.

Fig. 4.2). Such a flow-mod message will add the *encap* action to the entry it modifies (Line 8-10 in Alg. 2). If the predecessor switch is a packet-in switch (*sw1* in Fig. 4.2), the AC is carried by the flow-mod/packet-out or the packet-out message (Label 6 in Fig. 4.2). In this case, the *encap* action will apply to the packet directly, since the packet in these messages is the first data packet. (Line 2-5 in Alg. 2 and Line 2-4 in Alg. 3, respectively)

### 4.3.3 Addressing type-3 race – Gurantee one packet-in for each new flow

The type-3 race results from duplicated packet-in messages. As we can see that in Attendre each new flow only generates one packet-in. So the type-3 race is eliminated.

### 4.3.4 Handling multiple flow-mods to a single switch – Message Set

In Attendre, the switch removes the buffer action and lets the buffered packets to re-match the flow tables after the switch processes the command messages for the flow. In some cases, the controller application might send multiple messages to a switch for one flow. In order to release the buffered packets at the right time, each switch needs to know how many messages are sent to it. To address this issue, we introduce the notion of Message Set (MS). An MS is a collection of messages the controller sends to a switch in response to a packet-in message. An MS has an MS version for identification, and an MS size telling the switch the number of messages in that MS. Since the MS version is only used locally at each switch, the controller could assign the same version for the MSs resulting from processing packet-in message. The MS version and size are included in the ACs and the flow-mod and packet-out messages. With the MS version and size, the switch could remove the *buffer* entry of the flow and release the packets after the switch has processed all the messages in the MS.

(Line 13-16 in Alg. 2 and Line 6-9 in Alg. 3, respectively)

#### 4.3.5 Handling flow update – MS version

In a flow update, the controller might send a flow-mod to modify the action field of an existing entry in the flow table. When a following hop switch gets an AC, if an entry in the switch has the same match fields and priority value as the AC, the switch only replaces the action with the *buffer* action when the MS version of the AC is higher than that of the entry; if the MS versions are the same, the *buffer* action will be appended to the entry if the switch has not processed all the flow-mod and packet-out messages in the MS specified by the MS version of the AC; otherwise, the switch ignores the AC. (Line 4-14 in Alg. 1)

#### 4.3.6 Handling multiple flow tables

In OpenFlow version 1.3, a switch has multiple flow tables in a pipeline. For such a switch with multiple flow tables, we highlight the following Attendre features. 1) At the packet-in switch, the switch inserts the *buffer* entry for the mismatched packet in the flow table that the packet mismatches with. 2) At a following hop switch, the *buffer* entry should be inserted to the earliest flow table the controller modifies. 3) When the buffered packets are released, they starts matching the flow table where the *buffer* entry is inserted.

## Chapter 5

# Implementation of Attendre

### 5.1 Implementation in Open vSwitch

We implement Attendre in Open vSwitch v1.9.3, the most recent LTS version [11], with about 500 lines of code. Open vSwitch (OVS) is a software OpenFlow switch, which consists of a *kernel* module and an *userspace* module. The kernel module emulates the OpenFlow data plane for fast packet processing. It receives packets, matches the packets against the kernel flow table, and processes the packets according to the actions defined in the matched entries. If a packet cannot find a matched entry in the kernel flow table, the packet will be sent to the userspace via inter-process communication. The userspace module implements the OpenFlow protocol. In particular, the userspace module connects to the OpenFlow controller, and creates, sends, receives and processes OpenFlow messages, e.g., packet-in, packet-out, flow-mod, and barrier request and reply. The userspace module also configures and manages the kernel flow table according to the OpenFlow messages it receives.

The Attendre features are added to the OVS userspace module. No change to the OVS kernel module is needed. More specifically, the packet buffering, encapsulation and decapsulation and the *buffer* entry insertion and deletion are all done in the userspace module. Implementing Attendre in the OVS userspace module allows Attendre to run on many commercial hardware switch platforms from companies like Marvell, Broadcom, Cavium [14] and Pica8 [15] that are compatible with the OVS



userspace module. Implementing Attendre in the userspace module also allows it to use the switch CPU memory, which is typically implemented by cheap DRAM with large capacity, for packet buffering, instead of using the packet buffer on the switch ASIC chip which is usually implemented by expensive SRAM with very limited capacity.

## 5.2 Implementation in NOX

NOX is an OpenFlow controller platform [12]. We implement the changes involved in Attendre in NOX. In Attendre, the controller issues the packet-out and the flow-mod messages carrying the AC if it wants to insert a *buffer* entry in the next hop switch. We add APIs with which the application can create these new types of messages. The MS version and size are assigned to the command messages by the controller itself.

## 5.3 Implementation in NICE

NICE is a model checking based OpenFlow application verification tool. In NICE, the OpenFlow network components, such as data packet, command message, host, OpenFlow switch and NOX controller platform, are implemented as different models. We integrate Attendre into NICE by modifying these network component models accordingly.

## Chapter 6

### Evaluation

#### 6.1 Benefits to Model Checking Verification

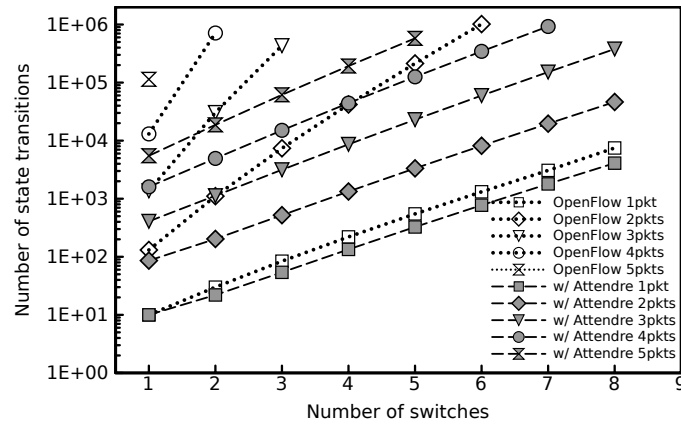
We use NICE to verify a few representative applications including “route flow”, mac learning switch, energy-efficient traffic engineering. We compare the number of state transitions in the verification and the verification execution time in OpenFlow and Attendre. The machine doing the verifications has a Quad-Core AMD Opteron 2393 SE CPU and 16GB of RAM.

##### 6.1.1 Route Flow

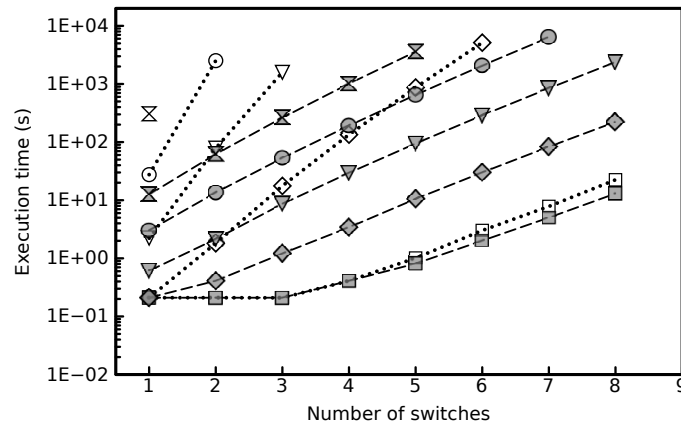
The verified network has a sender and a replier. The sender sends data packets and receives the same number of packets returned by the replier. We run the verification with different numbers of packets and different numbers of switches between the sender and the replier.

The number of state transitions and the execution time incurred during the model checking for different configurations (up to eight switches and one million state transitions) are shown in Fig. 6.1. Fig. 6.1(b) shares the same legend with Fig. 6.1(a). NICE uses a coarse time granularity of 0.2s to measure the execution time. So the execution time never falls below 0.2s.

Fig. 6.1 shows that in Attendre, the verification cost, i.e., the number of state transitions and the execution time, can be reduced by several orders of magnitude.



(a) Number of state transitions



(b) Execution time

Figure 6.1 : NICE verification results under various scenarios for OpenFlow and Attendre

For example, for the case having two packets and six hops between the sender and the replier (diamond markers), the cost are reduced by more than two orders of magnitude (about 200 times).

We can also glean the benefits of buffering the packets at switches from Fig. 6.1. The benefit of buffering the non-initial packets at the first switch is shown by the experiments with only one switch. For example, in the case with one switch and five packets (cross markers), the costs decrease by at least 20 times with Attendre. The

benefit of buffering the initial packet at the following switches can be seen in the experiments with only one packet (square markers). For example, if there are eight switches on the path, Attendre could save the verification costs by about 50%.

### 6.1.2 MAC Learning Switch

This application achieves the mac address learning protocol for individual switch. It learns the mac address connectivity by checking the source mac address of the packet in the packet-in message, and floods the packet if the destination address is unknown, otherwise, it installs a rule forwarding the packet to the learned port. The experiment uses a line topology and varies the number of packets sent by the sender and the number of hops between the sender and replier. The replier could move to another port on the switch at a random time. Since the MAC learning application only installs rules for one hop individually, there is only type 1 race in the network.

Table 6.1 shows the verification cost of this application under different configurations. For the cases with one switch, the verification cost reduces by about five times when there are four packets. When the number of switches increases, with the same number of packets, it has even more reduction. The verification cost for the cases with one packets keeps the same between Attendre and OpenFlow, since there is no race in these cases.

### 6.1.3 Energy-Efficient Traffic Engineering

SDN allows the controller application to dynamically select the forwarding path for a flow according to the traffic load on the switch and power down the unnecessary links to reduce the energy consumption. We verify REsPoNse [16], which pre-computes the “alwaysOn” and “onDemand” path for each flow. When the REsPoNse controller

Table 6.1 : Number of State Transitions and Execution Time in MAC Learning Application

transition /time(second)	1 switch		2 switches	
	Attendre	OpenFlow	Attendre	OpenFlow
1 pkt	74/0.2	74/0.2	188/0.5	188/0.5
2 pkts	1.8k/2.5	2.1k/2.3	9.2k/18	11k/20
3 pkts	16.5k/27	36.1k/54	133k/355	335k/847
4 pkts	101k/209	527k/1019	-	-

application receives a packet-in message, it selects the path for the flow according to the switch port statistics value, e.g., if the switch is under light traffic load it chooses the “alwaysOn” path, otherwise it chooses the “onDemond” path. We use the symbolic engine in NICE to generates different combination of the port statistics In the experiments, we vary the number of switches on the “alwaysOn” and “onDemond” path between the sender and repliers and the number of flows the sender sends. The sender only sends one packet, so only type 2 race happens in these cases.

Table 6.2 shows the verification costs. Since the controller application installs a forwarding path for the packets, the type 2 race in OpenFlow results in many redundant command messages. From Table 6.2, we find that Attendre has much lower verification cost than that in OpenFlow, since in all these cases, each flow only creates one packet-in message. FOr example, the verification cost can be reduced by 6 times when there are three switches on the “alwaysOn” path and 4 switches on the “onDemond” path. As the cost in verifying the route flow application, the benefit would be even more when the number of switches on the paths increases.

Table 6.2 : Number of State Transitions and Execution Time in Energy-Efficient Traffic Engineering Application

transition /time(second)	alwaysOn: 2 hops		alwaysOn: 3 hops	
	onDemand: 3 hops		onDemand: 4 hops	
	Attendre	OpenFlow	Attendre	OpenFlow
1 flow	2.1k/23	2.9k/32	8.3k/162	47.1k/1043
2 flows	27.1k/364	38.0k/497	110k/2538	642k/16473

## 6.2 Benefits to TCP connection setup

This experiment aims at showing that Attendre could reduce the time needed to establish TCP connections. We first describe our experiment setup and then analyze the experiment results.

### 6.2.1 Experiment Setup

We setup a network prototype via a container-based network emulator Mininet v2.1.0 [17], with Open vSwitch and the NOX [12] OpenFlow controller platform. In this network, a client establishes a set of TCP connections with a server, i.e., the client sends TCP SYN packets with different source port numbers and receives the corresponding TCP SYN-ACK packets replied from the server. The experiment is done on a single desktop computer, which has an AMD Phenom II X4 965 CPU (4 cores each at 3.4GHz) and 8GB of RAM.

The TCP connection establishment time can be affected by the number of hops between the client and the server, and the processing delay of the switches on the forwarding path. The experiments in this section focus on studying the influence of these two factors. To compare Attendre with OpenFlow and OpenFlow-B in various

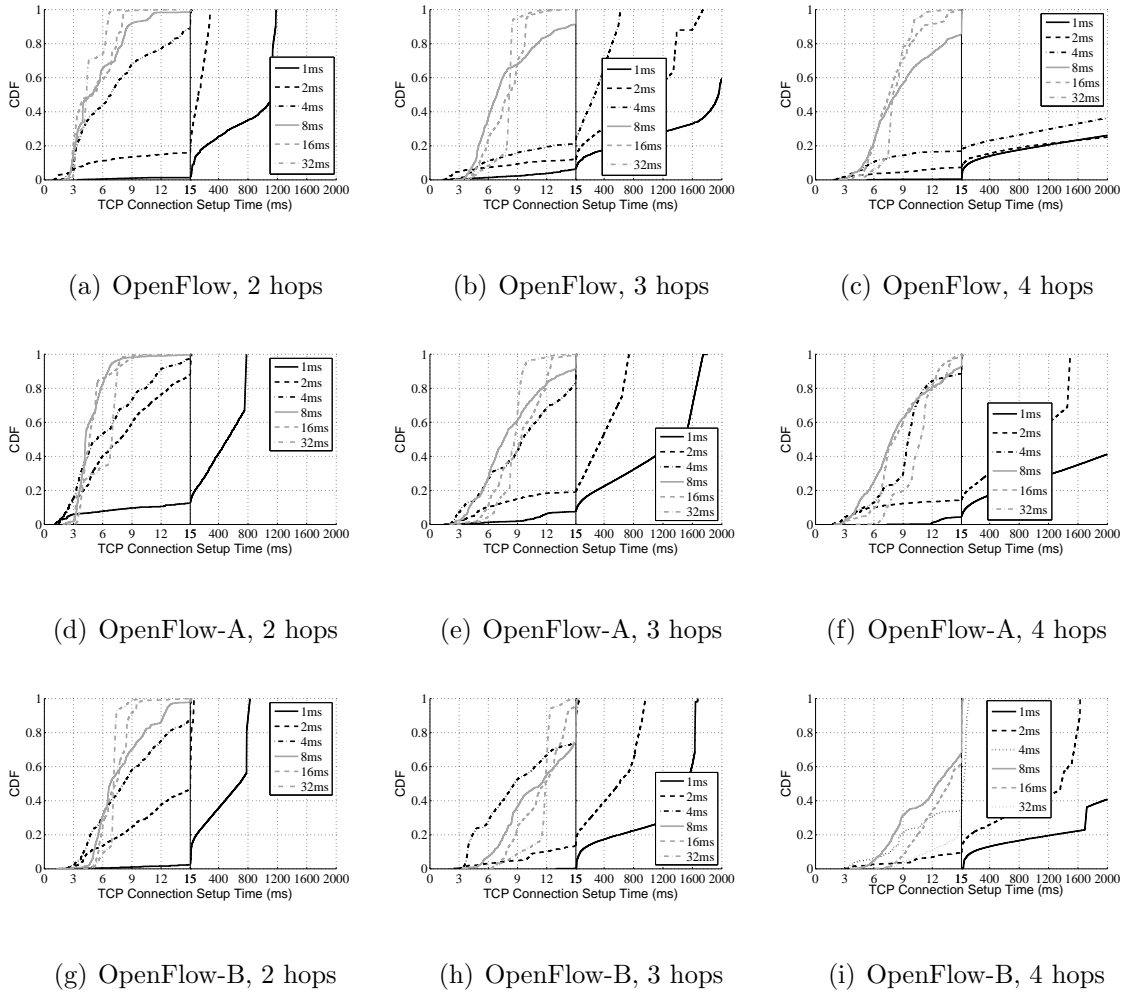


Figure 6.2 : CDF of TCP connection establishment time with different hops and flow inter-arrival time

situations, we vary the number of hops between the client and the server and the time interval between the flows, i.e., the flow inter-arrival time.

We measure the round trip time (RTT) of each pair of SYN and SYN-ACK messages so as to reflect the time needed to establish a TCP connection. The time needed by the third packet, TCP ACK, in the TCP 3-way handshake is ignored, since it consumes very little time comparing with the time needed by SYN and SYN-ACK. This is because ACK has the same match fields as SYN in the same connection, and the

flow table entry for them has already been installed in OVS kernel flow table when SYN traverses the network, so that ACK will not be passed to the OVS userspace.

The goal of this experiment is to analyze the TCP connection setup time for each pair of client and server in an SDN network where many clients make connections to many servers. However, from the network perspective, using only a single pair of client and server to emulate the traffic has no difference from using multiple clients and servers. This is because all the packets sent by the client and the server have different port numbers. By installing exact matching rules to the switches, the controller could differentiate these packets into different flows, as if they come from multiple clients and servers. The other reason is that having only one pair of client and server allows the flow arrival interval to be controlled more accurately, and also gets rid of the unnecessary overhead of the additional hosts.

The OpenFlow controller platform we use is NOX v0.9.1 [12]. NOX provides very small controller delay for request rates lower than 20000 requests per second [18], which is higher than the packet-in request rate in our experiment. The application running on the NOX controller is a layer 2 forwarding application that installs exact matching rules to switches on the forwarding path in response to each packet-in. This forwarding application is based on a network topology learning application running in parallel. We implement this forwarding application for OpenFlow, Attendre and OpenFlow using barrier messages respectively.

It is worthy to note that, if the controller is heavily loaded or has long processing delay, Attendre will show more benefits in packet forwarding delay. The influence of the controller processing delay will be discussed in Chapter 6.3.

Even though the experiments are done with software OpenFlow switch, we predict that the hardware switch running OVS would show the same trends and similar



results. The reason is that, as we discussed in Chapter 5, the difference between OVS switches with and without Attendre is in the userspace. This difference is exactly the same difference between the hardware switches running OVS as the control plane.

For the switch implementing Attendre into the hardware, the *buffer* and *encap* TCAM rule insertion takes about tens of memory movement [19]. Thus, the time cost by the hardware rule insertion is at the scale of nanosecond, which is comparable to software rule insertion and deletion time. We expect that Attendre should still show similar performance improvement with the hardware implementation.

### 6.2.2 Results

In the experiment, for each case, the client establishes a thousand TCP connections with the server. Each case has a constant flow arrival interval varying from 1ms to 32ms and the number of hops between the client and the server varies from two to four. The CDFs of the 1000 RTTs for all the cases are shown in Fig. 6.2. In general we can conclude that the more switches on the path, the higher the RTTs are; and the lower flow arrival interval, the higher the RTTs are. Besides these observations, we make the following observations by comparing Attendre with OpenFlow and OpenFlow-B.

**When the flow arrival interval is small (less than 16ms), Attendre has much lower RTTs than OpenFlow**

When the flow arrival interval is small, the switch in the network needs to send more packet-ins and process more flow-mods per unit time, which increases the delay of sending or processing the messages. In **OpenFlow**, this longer delay increases the chance that a packet arrives at a following switch where the corresponding flow-mod has not been processed. If a packet is sent to the controller by a following switch, its

RTT will increase a lot. Moreover, when the controller receives a packet-in from a following switch, it will again send a set of flow-mods back, which further increases the number of flow-mods that the switches need to process per unit time.

**Attendre** buffers the mismatching packet and forwards it after the corresponding flow-mod is processed. This avoids sending redundant packet-ins and reduces the number of flow-mods a switch needs to process to the minimum number. Although Attendre adds a little overhead to the flow-mod processing and packet-in sending (i.e., inserting and removing the wait entries, encapsulating and de-encapsulating the packets and updating the Attendre tables), the comprehensive effect of the number of commands processed makes the workload per unit time of Attendre much less than OpenFlow. As a consequence, Attendre exhibits much less delay than OpenFlow. For example, for the case with 3 hops and 4ms flow arrival interval, the 80th percentile of the RTT is reduced from 500ms (not shown in the figure) to 15ms.

It is worthy to note that in a data center network traffic study [20], the flow inter-arrival time is quite low. For example, in the top-of-rack switches in cloud data centers, more than 95% of the flow inter-arrival times are less than 10ms; and the percentage in university data centers is 60-85%. Thus, Attendre could speed up the TCP connection times for most flows in data center networks.

### **In all cases, Attendre has lower RTTs than OpenFlow-B**

In OpenFlow-B, the controller sends the flow-mod/packet-out back to the first switch after it receives all the barrier replies from the following switches. Thus, although OpenFlow-B generates the same the numbers of packet-ins and flow-mods as Attendre and Attendre needs a little bit more time to processing a packet-in and flow-mod, waiting for the barrier replies adds much more delay to the RTTs. Moreover, pro-

cessing the barrier request and reply at the switch further delays the processing of other messages.

### **When the flow arrival interval is high, Attendre has a little larger RTTs than OpenFlow**

For example, for the case of three switches and a flow arrival interval of 32ms, the 80th-percentile of the RTT in Attendre is only about 2ms larger than that in OpenFlow. This performance degradation in Attendre is due to the fact that Attendre adds a little overhead to the processing of OpenFlow messages. On the other hand, for OpenFlow as the flow arrival interval is high, in a specific period of time, the number of flow-mod messages a switch needs to process is low. Thus, in most cases, the following switches are able to finish processing the flow-mod messages before the corresponding packets arrive at the switch. As a consequence, the chance of sending packet-in to the controller by the following switches is quite low.

## **6.3 Benefits to Forwarding Delay when Controller and Network Delays Vary**

This section addresses the issue about how the controller delay and the network delay between the controller and the switches affect the forwarding time of a packet in SDN network. We first derive the forwarding delay formulae for an unicast packet traversing  $N$  hops for three different cases, i.e., OpenFlow, OpenFlow-B and OpenFlow-Attendre. Then, based on these formulae, we numerically analyze the forwarding delay and show that Attendre could reduce the forwarding delay in most scenarios.

Table 6.3 : Delay Formulae Notation

Notation	Meaning
$N$	Number of <i>switches</i> on the path
$T_{case}^N$	Forwarding delay of a specific <i>case</i> with $N$ switches on the path (3 cases followed)
$nb$	OpenFlow not using <i>no barrier</i> message
$b$	OpenFlow using <i>barrier</i> message
$a$	OpenFlow with <i>Attendre</i> mechanism
$L_k$	Delay from $node_k$ pushing the packet in its output queue to $node_{k+1}$ beginning to process the packet
$S$	Delay from $host_0$ pushing the packet in its output queue to <i>controller</i> beginning to process the packet
$\tau^k$	Round-trip TCP delay between $switch_k$ and <i>controller</i>
$g^k$	Propagation delay between $node_k$ and $node_{k+1}$
$t^k$	Transmission delay from $node_k$ and $node_{k+1}$
$p_{sp}^k$	Processing delay for $switch_k$ processing <i>packet</i> (from matching the packet against flow table to executing the action to the packet)
$p_{sc}^k$	Processing delay for $switch_k$ processing <i>command</i>
$p_c^k$	Processing delay for <i>controller</i> processing packet-in message from $switch_k$
$q_i^k$	Queueing delay of the <i>incoming</i> packet for $switch_k$
$q_o^k$	Queueing delay of the <i>outgoing</i> packet for $switch_k$

### 6.3.1 Delay Model in OpenFlow Network

The delay in traditional network architecture is composed of four components including propagation ( $g^k$ ), transmission ( $t^k$ ), queueing ( $q_i^k/q_o^k$ ) and processing delay. These delay components are also applicable to the delay computation of the OpenFlow network except for the processing delay containing three separate ingredients, which are  $p_{sp}^k$ ,  $p_{sc}^k$  and  $p_c^k$ . Besides that, the packets in OpenFlow could experience several TCP

round trip time ( $RTT = \tau^k$ ) between  $switch_k$  and  $controller$ , since the packet will be sent to the controller if the packet doesn't match the flow table.

### 6.3.2 Delay Formulae

The notations used in this section are illustrated in Table 6.3. In the table, *node* means the *switches* or *hosts* in the network. Indices 0 and  $N + 1$  are for the sender and the receiver respectively, and indices from 1 to  $N$  are for the switches. To make the formulae easier to understand, we also define some frequently used variables in (6.1) and (6.2) whose definitions are in Table 6.3.

$$L_k \equiv q_o^k + t^k + g^k + q_i^{k+1} \quad 0 \leq k \leq N \quad (6.1)$$

$$S \equiv L_0 + p_{s_p}^1 + \frac{\tau^1}{2} + p_c^1 \quad (6.2)$$

The formula for **OpenFlow** is shown in (6.3). Since the matched entry is absent on the switches, in all of the three cases, the packet will be sent to *controller* by  $switch_1$ . Thus, if there is one switch on the path ( $N = 1$ ), the delay should include  $S$ , the delay of processing the packet-in, the TCP delay of sending the flow-mod/packet-out back, the delay of processing the message and  $L_1$ . If there are more than one switch ( $N \geq 2$ ) on the path, the delay with  $N$  switches contains the delay with  $N - 1$  switches, the processing delay of the packet at  $switch_N$  and  $L_N$ . Besides these, if the packet is processed by  $switch_N$  before it process the flow-mod, another RTT between the switch and the controller and corresponding processing delays should be taken into account, which is the last term in (6.3).  $I(b)$  is an indicator function. The value

of the function is 1 when  $b$  is *true* and 0 when  $b$  is *false*.

$$T_{nb}^N = \begin{cases} S + p_{s_c}^1 + \frac{\tau^1}{2} + p_{s_p}^1 + L_1 & N = 1 \\ T_{nb}^{N-1} + p_{s_p}^N + L_N + (\tau^N + p_c^N + p_{s_c}^N + p_{s_p}^N) \\ \quad \times I((T_{nb}^{N-1} - S) < (\frac{\tau^N}{2} + p_{s_c}^N)) & N \geq 2 \end{cases} \quad (6.3)$$

The formula for **OpenFlow-B** is shown in (6.4). The controller waiting time is the last term. The first two summations are for the delay on the forwarding path. The following four terms are the round trip time for packet-in and flow-mod and the corresponding processing delays.

$$T_b^N = \sum_{k=0}^N L_k + \sum_{k=1}^N p_{s_p}^k + p_{s_p}^1 + \tau^1 + p_c^1 + p_{s_c}^1 + \max_{2 \leq j \leq N} (\tau^j + p_{s_c}^j) \quad (6.4)$$

The formula for **OpenFlow-Attendre** is shown in (6.5). For  $N = 1$ , the delay is the same as the other two cases. For  $N \geq 2$ , the first term in the max function is the delay from the time when *controller* sends the flow-mod to *switch*<sub>1</sub> to the time *switch* <sub>$k$</sub>  finishes processing the packet encapsulated with the AC. The second term in the max function is the delay from the time when *controller* sends the flow-mod to *switch* <sub>$k$</sub>  to the time when *switch* <sub>$k$</sub>  finishes processing that flow-mod. The maximum value of these two terms plus  $S$  is the delay from the beginning to the time *switch* <sub>$k$</sub>

starting processing the packet.

$$T_a^N = \begin{cases} S + \frac{\tau^1}{2} + p_{s_c}^1 + p_{s_p}^1 + L_1 & N = 1 \\ p_{s_p}^N + L_N + S \\ \quad + \max((T_a^{N-1} - S), (\frac{\tau^N}{2} + p_{s_c}^N)) & N \geq 2 \end{cases} \quad (6.5)$$

### 6.3.3 Experiment Setup

We numerically analyze the one-way forwarding time needed by a single initial packet. The word *initial* here means that the forwarding entries matching this packet are not in the flow tables of the switches, so that the controller will install the flow table entries after it receives a packet-in message. This packet could be TCP SYN packet in reality. We deliberately choose the delay values for different network delay components, such as transmission, propagation, queuing and processing, to imitate data center network scenarios. In addition, we also carefully choose the values for the TCP link delay between the switches and the controller.

The **number of switches** on the forwarding path is 6, since 6 hops can connect most servers in a medium size data center, like the 2000-node Hadoop cluster in Facebook for spam detection and ad optimization [21], with commodity 48 port switches [22] organized in a fat tree topology [23] [24]. The **transmission delays** set here are for a TCP SYN packet, since Microsoft reported that in their data center, 99.91% of traffic is TCP [22]. So the transmission delays are  $3.04\mu s$  and  $0.576\mu s$  for encapsulated and non-encapsulated SYN respectively, by assuming a 1Gbps port speed. The per hop **propagation delay** in a data center is trivial comparing to other delays, since the signal propagating through a 50m copper cable only takes  $167ns$ .

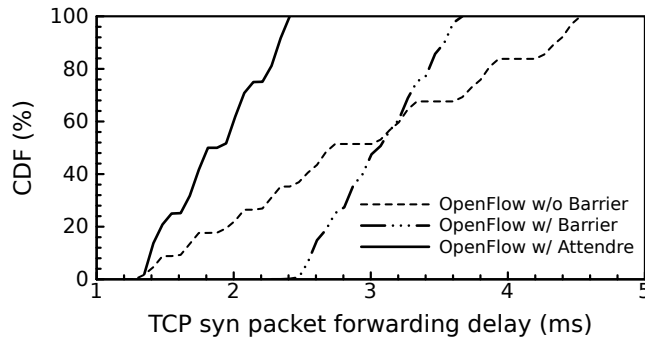


Figure 6.3 : Delay Analysis

We use this number as an approximation for the propagation delay. The **queuing delay** we use is  $3.4\mu s$ , which is half of the transmission delay for a packet with the average packet size, 850 bytes [25]. The **data packet processing delay** at the OpenFlow switch is no more than  $5\mu s$  [26]. The **command message processing delay** at the switch is about  $1ms$  [27].

We keep those aforementioned delays fixed through our computation. Then, we vary the **controller packet-in message processing delay** and the **TCP link delays** between the switches and the controller. To imitate the different controller workloads which result in various controller processing delays, we vary this value from  $0.1ms$  to  $1ms$  [18] by taking four numbers uniformly distributed in this interval. The physical position of the controller could be varied. So the round trip TCP delays between switches and the controller take on six possible values uniformly distributed between  $50\mu s$  and  $250\mu s$ , since the TCP delay in a data center is no more than  $250\mu s$  [22].



### 6.3.4 Numerical Analysis Results

The numerical analysis is performed on all combinations of these various delay values and applies each combination to three different OpenFlow mechanisms, OpenFlow, OpenFlow-B and OpenFlow-Attendre. Fig. 6.3 shows the results based on the parameters we choose and the formulae. Even though Attendre encapsulates AC to packets and buffer packets, the initial packet forwarding delay in OpenFlow-Attendre is smaller than other cases. This is because with Attendre, packets are buffered at the switches to avoid unnecessary controller processing delay, so that the packets could be processed at a switch just after the corresponding flow-mod messages are processed by the switch.

More specifically, in OpenFlow-Attendre, the initial packet forwarding delays are less than OpenFlow-B by about 1 ms. The worst case of OpenFlow-Attendre is actually better than the best case of OpenFlow-Barrier. Comparing OpenFlow-Attendre with OpenFlow, the worst case of OpenFlow-Attendre is better than 60% of the cases of the latter. This delay reduction could be very significant to certain deployment scenarios, e.g., financial institutions, where end-to-end application performance is sensitive to delays of millisecond or less. Anecdotally, the economic loss in a stock trading system could be \$4 million in revenues per millisecond of extra delay [28].

## Chapter 7

### Related Work

#### 7.1 Previous work related to SDN verification

The issue of SDN verification is being explored actively in academia. Kinetic [29] proposed a way to consistently update the forwarding policies in OpenFlow network, i.e., the packets or flows are processed by each hop with the same version of a policy. With the consistency guarantee, Kinetic simplifies the dynamic flow table configuration verification to a low cost static verification. Comparing with Kinetic, Attendre has three major differences. **First**, the verification approach that Attendre benefits to, model checking, is a more complete and rigorous verification approach than flow table configuration verification. The latter one only verifies reachability problems like is there a forwarding loop, does a packet go through a specific path or can host A talk to host B. In contrast, model checking could verify more sophisticated invariant checking questions like are the packets from a specific flow sent to the controller more than once or is two packets reordered by the network. **Secondly**, in Kinetic, the consistency guarantee is achieved by tagging the flow table entry with a version number and tagging the version number to the packets at the first switch so that the packets could only match the entry with the tagged version number. However, when the controller installs a forwarding policy for a new flow, Kinetic changes the ingress port configuration of the first switch only after the controller knows that all of the flow-mod messages to the other switches have been installed. This way of performing

network update is essentially identical to using a barrier message as described in Section 3.3. We already show that the OpenFlow-B like method increases the forwarding delay of the packet. **Thirdly**, although Kinetic eliminates the ill effects of the type 2 race, it does not consider type 1 and 3 races. Thus, in Kinetic, the first switch could still send redundant packet-in messages to the controller.

Another SDN verification method is run-time verification. VeriFlow [30] and NetPlumber [31] belong to this category. They check the run time behaviors of applications in an operating network against packet processing invariants, and block the behaviors that violate the invariants. However, these run time tools cannot be used to verify applications before they are deployed. Thus, a potentially buggy application may be deployed live, i.e., bugs are only caught and blocked after the fact, and the network can only continue to function in a degraded mode. For these reasons, the focus of our analysis of OpenFlow has been on its suitability for model checking based application verification, an approach that is being actively developed by the NICE project [10], and is the most powerful among the existing alternatives.

Guha *et al.* proposes a machine-verified SDN controller [32] for NetCore [33], a domain specific language for SDN controller programming. This controller platform could verify if the NetCore compiler generates the correct flow table rules given a controller application written in NetCore. However, the controller does not check the correctness of the application.

## 7.2 Previous work aiming at reducing the first packet forwarding delay and control plane workload in SDN

The aspect of reducing the first packet forwarding delay and the control plane workload are delved into actively. Different from OpenFlow, DIFANE [34], another enterprise network architecture, reduces the control plane workload by pre-computing forwarding rules and distributing them to the data plane of some “authority” switches. Instead of sending the mismatching packet to the controller, the first switch will send it to an authority switch, which directly forwards the packet to its destination and installs rules back to the first switch. Thus, the first packet forwarding delay can be reduced by always keeping packets in the data plane. However, DIFANE restricts that the forwarding rule should be pre-computed and installed in the authority switch, which makes it reacts to dynamic network changes slowly. DevoFlow [35] is another modification to OpenFlow protocol that can reduce the first packet forwarding delay and the control plane workload. In DevoFlow, switches will make local forwarding decision for small flows, whereas only the large flow are sent to the centralized controller, which means that DevoFlow limits the visibility of the controller to only significant flows.

Both DIFANE and DevoFlow point out that a single physical controller is the bottleneck of SDN, and successfully reduce the control plane workload and the first packet forwarding delay. However, they sacrifice the ability for the controller to inspect all flows that arrive at the network. In contrast, Attendre keeps this ability unchanged. Moreover, none of them address the issue with SDN application verification efficiency.

Kandoo [36] investigates another inefficiency in SDN control plane, i.e., processing

the network events does not always use network-wide state. To relief the workload on the single centralized controller, Kandoo divides the control plane into two layers. The controller in the top layer runs the network application using the knowledge of network-wide state. While the bottom layer of the control plane contains a group of local controllers that run the applications only use local information that can be retrieved from a single switch. However, Attendre solves the problem of redundant messages in the control plane, which is orthogonal to the inefficiency that Kandoo addresses. This difference also makes Attendre compatible with Kandoo. With Attendre, both the bottom layer and the top layer controllers in Kandoo should expect less workload.

## Chapter 8

### Conclusion

We have identified three types of race condition present in an OpenFlow network, and presented a mechanism called *Attendre* that mitigates the ill effects of these race conditions. As our results indicate, the benefit of *Attendre* is potentially very large. More broadly, this work emphasizes the importance of considering race condition in SDN protocol design and shows that by taking the consequence of race conditions as a first-class protocol design concern, it is possible to gain a deeper understanding of the subtle behaviors of existing protocols, and to arrive at improved protocol designs.

## Bibliography

- [1] S. Jain, A. Kumar, S. Mandal, J. Ong, L. Poutievski, A. Singh, S. Venkata, J. Wanderer, J. Zhou, M. Zhu, J. Zolla, U. Holzle, S. Stuart, and A. Vahdat, “B4: experience with a globally-deployed software defined wan,” in *SIGCOMM'13*.
- [2] T. Hinrichs, J. Mitchell, N. Gude, S. Shenker, and M. Casado, “Expressing and enforcing flow-based network security policies,” tech. rep., University of Chicago, 2008.
- [3] S. A. Mehdi, J. Khalid, and S. A. Khayam, “Revisiting traffic anomaly detection using software defined networking,” in *RAID'11*.
- [4] S. Shin, P. Porras, V. Yegneswaran, M. Fong, G. Gu, and M. Tyson, “Fresco: Modular composable security services for software-defined networks,” in *NDSS'13*.
- [5] P. Porras, S. Shin, V. Yegneswaran, M. Fong, M. Tyson, and G. Gu, “A security enforcement kernel for openflow networks,” in *HotSDN '12*.
- [6] M. Casado, M. J. Freedman, J. Pettit, J. Luo, N. Gude, N. McKeown, and S. Shenker, “Rethinking enterprise network control,” *IEEE/ACM Trans. Netw.*, vol. 17, pp. 1270–1283, Aug. 2009.
- [7] A. K. Nayak, A. Reimers, N. Feamster, and R. Clark, “Resonance: Dynamic access control for enterprise networks,” in *WREN '09*.

- [8] M. Reitblatt, M. Canini, A. Guha, and N. Foster, “Fattire: Declarative fault tolerance for software-defined networks,” in *HotSDN '13*.
- [9] P. Gill, N. Jain, and N. Nagappan, “Understanding network failures in data centers: measurement, analysis, and implications,” in *SIGCOMM '11*.
- [10] M. Canini, D. Venzano, P. Perešini, D. Kostić, and J. Rexford, “A nice way to test openflow applications,” in *NSDI'12*.
- [11] “Open vswitch-1.9.3.” <http://openvswitch.org/download/>, 2013.
- [12] N. Gude, T. Koponen, J. Pettit, B. Pfaff, M. Casado, N. McKeown, and S. Shenker, “Nox: Towards an operating system for networks,” *SIGCOMM Comput. Commun. Rev.*, vol. 38, pp. 105–110, July 2008.
- [13] M. Casado, M. J. Freedman, J. Pettit, J. Luo, N. McKeown, and S. Shenker, “Ethane: Taking control of the enterprise,” in *SIGCOMM '07*.
- [14] “xflowresearch.” <http://www.xflowresearch.com/>.
- [15] “Pica8.” <http://www.pica8.com/open-technology/open-vswitch.php>.
- [16] N. Vasić, P. Bhurat, D. Novaković, M. Canini, S. Shekhar, and D. Kostić, “Identifying and using energy-critical paths,” in *CoNEXT '11*.
- [17] N. Handigol, B. Heller, V. Jeyakumar, B. Lantz, and N. McKeown, “Reproducible network experiments using container-based emulation,” in *CoNEXT '12*.
- [18] Z. Cai, A. L. Cox, and T. S. E. Ng, “Maestro: A system for scalable openflow control,” Tech. Rep. TR10-08, Rice University, Dec. 2010.



- [19] D. Shah and P. Gupta, “Fast updating algorithms for tcams,” *IEEE Micro*, vol. 21, no. 1, pp. 36–47, 2001.
- [20] T. Benson, A. Akella, and D. A. Maltz, “Network traffic characteristics of data centers in the wild,” in *IMC '10*.
- [21] B. Hindman, A. Konwinski, M. Zaharia, A. Ghodsi, A. D. Joseph, R. Katz, S. Shenker, and I. Stoica, “Mesos: a platform for fine-grained resource sharing in the data center,” in *NSDI'11*.
- [22] M. Alizadeh, A. Greenberg, D. A. Maltz, J. Padhye, P. Patel, B. Prabhakar, S. Sengupta, and M. Sridharan, “Data center tcp (dctcp),” in *SIGCOMM'10*.
- [23] M. Al-Fares, A. Loukissas, and A. Vahdat, “A scalable, commodity data center network architecture,” in *SIGCOMM'08*.
- [24] A. Singla, C.-Y. Hong, L. Popa, and P. B. Godfrey, “Jellyfish: Networking data centers randomly,” in *NSDI'12*.
- [25] T. Benson, A. Anand, A. Akella, and M. Zhang, “Understanding data center traffic characteristics,” *SIGCOMM Comput. Commun. Rev.*, vol. 40, pp. 92–99, Jan. 2010.
- [26] C. Rotsos, N. Sarrar, S. Uhlig, R. Sherwood, and A. W. Moore, “Oflops: An open framework for openflow switch evaluation,” in *PAM'12*.
- [27] M. Jarschel, S. Oechsner, D. Schlosser, R. Pries, S. Goll, and P. Tran-Gia, “Modeling and performance evaluation of an openflow architecture,” in *ITC'11*.
- [28] cPacket Networkds, *Pragmatic Network Latency Engineering Fundamental Facts and Analysis*, 2009.

- [29] M. Reitblatt, N. Foster, J. Rexford, C. Schlesinger, and D. Walker, “Abstractions for network update,” in *SIGCOMM’12*.
- [30] A. Khurshid, X. Zou, W. Zhou, M. Caesar, and P. B. Godfrey, “Veriflow: Verifying network-wide invariants in real time,” in *NSDI’13*.
- [31] P. Kazemian, M. Chang, H. Zeng, G. Varghese, N. McKeown, and S. Whyte, “Real time network policy checking using header space analysis,” in *NSDI’13*.
- [32] A. Guha, M. Reitblatt, and N. Foster, “Machine-verified network controllers,” in *PLDI ’13*.
- [33] C. Monsanto, N. Foster, R. Harrison, and D. Walker, “A compiler and run-time system for network programming languages,” in *POPL ’12*.
- [34] M. Yu, J. Rexford, M. J. Freedman, and J. Wang, “Scalable flow-based networking with difane,” in *SIGCOMM’10*.
- [35] A. R. Curtis, J. C. Mogul, J. Tourrilhes, P. Yalagandula, P. Sharma, and S. Banerjee, “Devoflow: scaling flow management for high-performance networks,” in *SIGCOMM ’11*.
- [36] S. Hassas Yeganeh and Y. Ganjali, “Kandoo: A framework for efficient and scalable offloading of control applications,” in *HotSDN ’12*.