

TIME4: Time for SDN

Technical Report[†], February 2016

Tal Mizrahi, Yoram Moses*

Technion — Israel Institute of Technology

Email: {dew@tx, moses@ee}.technion.ac.il

Abstract—With the rise of Software Defined Networks (SDN), there is growing interest in dynamic and centralized traffic engineering, where decisions about forwarding paths are taken dynamically from a network-wide perspective. Frequent path re-configuration can significantly improve the network performance, but should be handled with care, so as to minimize disruptions that may occur during network updates.

In this paper we introduce TIME4, an approach that uses accurate time to coordinate network updates. TIME4 is a powerful tool in software-defined environments, that can be used for various network update scenarios. Specifically, we characterize a set of update scenarios called *flow swaps*, for which TIME4 is the optimal update approach, yielding less packet loss than existing update approaches. We define the *lossless flow allocation problem*, and formally show that in environments with frequent path allocation, scenarios that require simultaneous changes at multiple network devices are inevitable.

We present the design, implementation, and evaluation of a TIME4-enabled OpenFlow prototype. The prototype is publicly available as open source. Our work includes an extension to the OpenFlow protocol that has been adopted by the Open Networking Foundation (ONF), and is now included in OpenFlow 1.5. Our experimental results show the significant advantages of TIME4 compared to other network update approaches, and demonstrate an SDN use case that is infeasible without TIME4.

Time is what keeps everything from happening at once
– Ray Cummings

I. INTRODUCTION

A. It's About Time

The use of synchronized clocks was first introduced in the 19th century by the Great Western Railway company in Great Britain. Clock synchronization has significantly evolved since then, and is now a mature technology that is being used by various different applications, from mobile backhaul networks [3] to distributed databases [4].

The Precision Time Protocol (PTP), defined in the IEEE 1588 standard [5], can synchronize clocks to a very high degree of accuracy, typically on the order of 1 microsecond [3], [6], [7]. PTP is a common and affordable feature in commodity switches. Notably, 9 out of the 13 SDN-capable switch silicons listed in the Open Networking Foundation (ONF) SDN Product Directory [8] have native IEEE 1588 support [9]–[17].

In this work we introduce TIME4, a **generic** tool for using time in SDN. One of the products of this work is a new feature that enables timed updates in OpenFlow, and has been incorporated in OpenFlow 1.5. Furthermore, we present a class of update scenarios in which the use of accurate time is provably optimal, while existing update methods are sub-optimal.

B. The Challenge of Dynamic Traffic Engineering in SDN

Defining network routes dynamically, based on a complete view of the network, can significantly improve the network performance compared to the use of distributed routing protocols. SDN and OpenFlow [18], [19] have been leading trends in this context, but several other ongoing efforts offer similar concepts. The Interface to the Routing System (I2RS) working group [20], and the Forwarding and Control Element Separation (ForCES) working group [21] are two examples of such ongoing efforts in the Internet Engineering Task Force (IETF).

Centralized network updates, whether they are related to network topology, security policy, or other configuration attributes, often involve multiple network devices. Hence, updates must be performed in a way that strives to minimize temporary anomalies such as traffic loops, congestion, or disruptions, which may occur during transient states where the network has been partially updated.

While SDN was originally considered in the context of campus networks [18] and data centers [22], it is now also being considered for Wide Area Networks (WANs) [23], [24], carrier networks, and mobile backhaul networks [25].

WAN and carrier-grade networks require a very low packet loss rate. Carrier-grade performance is often associated with the term *five nines*, representing an availability of 99.999%. Mobile backhaul networks require a Frame Loss Ratio (FLR) of no more than 10^{-4} for voice and video traffic, and no more than 10^{-3} for lower priority traffic [26]. Other types of carrier network applications, such as storage and financial trading require even lower loss rates [27], on the order of 10^{-5} .

Several recent works have explored the realm of dynamic path reconfiguration, with frequent updates on the order of minutes [23], [24], [28], enabled by SDN. Interestingly, for voice and video traffic, a frame loss ratio of up to 10^{-4} implies that service must not be disrupted for more than 6 milliseconds per minute. Hence, if path updates occur on a per-minute basis,

[†]This report is an extended version of [1], which was accepted to IEEE INFOCOM '16, San Francisco, April 2016. A preliminary version of this report was published in arXiv [2] in May, 2015.

*Yoram Moses is the Israel Pollak academic chair at Technion.

then transient disruptions must be limited to a short period of no more than a few milliseconds.

C. Timed Network Updates

We explore the use of *accurate time* as a tool for performing coordinated network updates in a way that minimizes packet loss. Software-managed network can significantly benefit from using time for *coordinating* network-wide orchestration, and for enforcing a given *order* of events. We introduce TIME4, which is an update approach that performs multiple changes at different switches at the same time.

Example 1. Fig. 1 illustrates a flow swapping scenario. In this scenario, the forwarding paths of two flows, f_1 and f_2 , need to be reconfigured, as illustrated in the figure. It is assumed that all links in the network have an identical capacity of 1 unit, and that both f_1 and f_2 require a bandwidth of 1 unit. In the presence of accurate clocks, by scheduling S_1 and S_3 to update their paths at the same time, there is no congestion during the update procedure, and the reconfiguration is smooth. As clocks will typically be reasonably well synchronized, albeit not perfectly synchronized, such a scheme will result in a very short period of congestion.

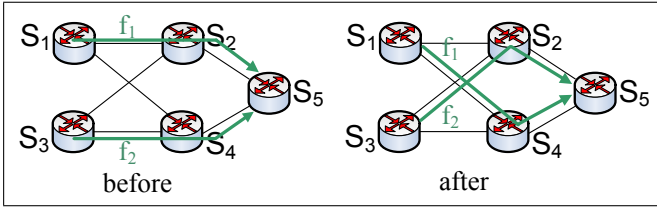


Fig. 1: Flow Swapping—Flows need to convert from the “before” configuration to the “after”.

In this paper we show that in a dynamic environment, where flows are frequently added, removed or rerouted, flow swaps are inevitable. A notable example of the importance of flow swaps is a recently published work by Fox Networks [29], in which accurately timed flow swaps are essential in the context of video switching.

One of our key results is that simultaneous updates are the optimal approach in scenarios such as Example 1, whereas other update approaches may yield considerable packet loss, or incur higher resource overhead. Note that such packet loss can be reduced either by increasing the capacity of the communication links, or by increasing the buffer memories in the switches. We show that for a given amount of resources, TIME4 yields lower packet loss than other approaches.

Accuracy is a key requirement in TIME4; since updates cannot be applied at the exact same instant at all switches, they are performed within a short time interval called the *scheduling error*. The experiments we present in Section IV show that the scheduling error in **software** switches is on the order of 1 millisecond. The TCAM-based **hardware** solution of [30] can execute scheduled events in existing switches with an accuracy **on the order of 1 microsecond**.

Accurate time is a powerful abstraction for SDN programmers, not only for flow swaps, but also for **timed consistent updates**, as discussed by [31].

D. Related Work

Time and synchronized clocks have been used in various distributed applications, from mobile backhaul networks [3] to distributed databases [4]. Time-of-day routing [32] routes traffic to different destinations based on the time-of-day. Path calendaring [33] can be used to configure network paths based on scheduled or foreseen traffic changes. The two latter examples are typically performed at a low rate and do not place demanding requirements on accuracy.

Various network update approaches have been analyzed in the literature. A common approach is to use a sequence of configuration commands [28], [34]–[36], whereby the **order** of execution guarantees that no anomalies are caused in intermediate states of the procedure. However, as observed by [28], in some update scenarios, known as **deadlocks**, there is no order that guarantees a consistent transition. **Two-phase** updates [37] use configuration version tags to guarantee consistency during updates. However, as per [37], *two-phase* updates cannot guarantee congestion freedom, and are therefore not effective in flow swap scenarios, such as Fig. 1. Hence, in flow swap scenarios the *order* approach and the *two-phase* approach produce the same result as the simple-minded approach, in which the controller sends the update commands as close as possible to instantaneously, and hopes for the best.

In this paper we present TIME4, an update approach that is most effective in flow swaps and other deadlock [28] scenarios, such as Fig. 1. We refer to update approaches that do not use time as **untimed** update approaches.

In SWAN [23], the authors suggest that reserving unused *scratch* capacity of 10-30% on every link can allow congestion-free updates in most scenarios. The B4 [24] approach prevents packet loss during path updates by temporarily reducing the bandwidth of some or all of the flows. Our approach does not require scratch capacity, and does not reduce the bandwidth of flows during network updates. Furthermore, in this paper we show that variants of SWAN and B4 that make use of TIME4 can perform better than the original versions.

A recently published work by Fox Networks [29] shows that accurately timed path updates are essential for video swapping. We analyze this use case further in Section IV.

Rearrangeably non-blocking topologies (e.g., [38]) allow new traffic flows to be added to the network by rearranging existing flows. The analysis of flow swaps presented in this paper emphasizes the requirement to perform *simultaneous* reroutes during the rearrangement procedure, an aspect which has not been previously studied.

Preliminary work-in-progress versions of the current paper introduced the concept of using time in SDN [39] and the flow swapping scenario [40]. The use of time for *consistent* updates was discussed in [31]. TimeFlip [30] presented a practical method of implementing timed updates. The current work is the first to present a generic protocol for performing timed updates in SDN, and the first to analyze *flow swaps*, a natural

application in which timed updates are the optimal update approach.

E. Contributions

The main contributions of this paper are as follows:

- We consider a class of network update scenarios called *flow swaps*, and show that simultaneous updates using synchronized clocks are provably the optimal approach of implementing them. In contrast, existing approaches for consistent updates (e.g., [28], [37]) are not applicable to flow swaps, and other update approaches such as SWAN [23] and B4 [24] can perform flow swaps, but at the expense of increased resource overhead.
- We use game-theoretic analysis to show that flow swaps are inevitable in the dynamic nature of SDN.
- We present the design, implementation and evaluation of a prototype that performs timed updates in OpenFlow.
- Our work includes an extension to the OpenFlow protocol that has been approved by the ONF and integrated into OpenFlow 1.5 [41], and into the OpenFlow 1.3.x extension package [42]. The source code of our prototype is publicly available [43].
- We present experimental results that demonstrate the advantage of timed updates over existing approaches. Moreover, we show that existing update approaches (SWAN and B4) can be improved by using accurate time.
- Our experiments include an emulation of an SDN-controlled video swapping scenario, a real-life use case that has been shown [29] to be infeasible with previous versions of OpenFlow, which did not include our time extension.

II. THE LOSSLESS FLOW ALLOCATION (LFA) PROBLEM

A. Inevitable Flow Swaps

Fig. 1 presents a scenario in which it is necessary to *swap* two flows, i.e., to update two switches at the same time. In this section we discuss the inevitability of flow swaps; we show that there does not exist a controller routing strategy that avoids the need for flow swaps.

Our analysis is based on representing the flow-swap problem as an instance of an unsplittable flow problem, as illustrated in Fig. 2b. The topology of the graph in Fig. 2b models the traffic behavior to a given destination in common multi-rooted network topologies such as fat-tree and Clos (Fig. 2a).

The unsplittable flow problem [44] has been thoroughly discussed in the literature; given a directed graph, a source node s , a destination node d , and a set of flow demands (commodities) between s and d , the goal is to maximize the traffic rate from the source to the destination. In this paper we define a *game* between two players: a **source**¹ that generates traffic flows (commodities) and a **controller** that reconfigures the network forwarding rules in a way that allows the network to forward all traffic generated by the source without packet losses.

¹The source player does not represent a malicious attacker; it is an ‘adversary’, representing the worst-case scenario.

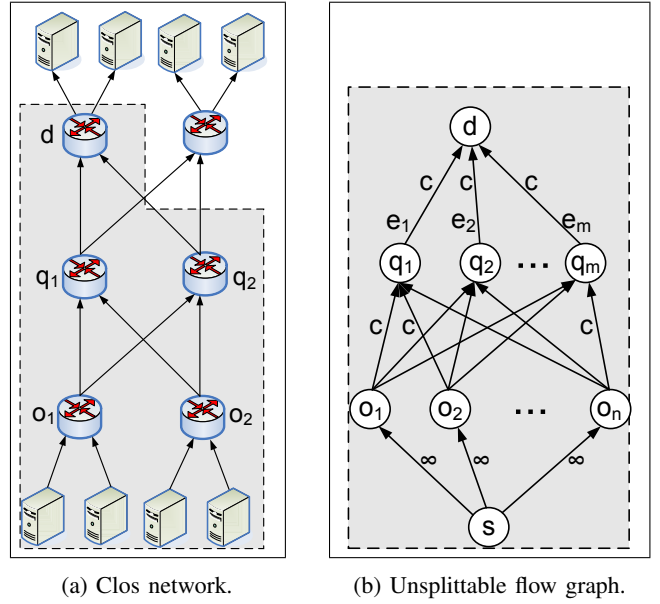


Fig. 2: Modeling a Clos topology as an unsplittable flow graph.

Our main argument, phrased in Theorem 1, is that the source has a strategy that **forces** the controller to perform a flow swap, i.e., to reconfigure the path of two or more flows at the same time. Thus, a scenario in which multiple flows must be updated at the **same time** is inevitable, implying the importance of timed updates.

Moreover, we show that the controller can be forced to invoke n individual commands that should optimally be performed at the same time. Update approaches that do not use time, also known as **untimed** approaches, cause the updates to be performed over a long period of time, potentially resulting in slow and possibly erratic response times and significant packet loss. Timed coordination allows us to perform the n updates within a short time interval that depends on the scheduling error.

Although our analysis focuses on the topology of Fig. 2b, it can be shown that the results are applicable to other topologies as well, where the source can force the controller to perform a swap over the edges of the min-cut of the graph.

B. Model and Definitions

We now introduce the *lossless flow allocation (LFA)* problem; it is not presented as an optimization problem, but rather as a game between two players: a **source** and a **controller**. As the source adds or removes flows (commodities), the controller reconfigures the forwarding rules so as to guarantee that all flows are forwarded without packet loss. **The controller’s goal** is to find a forwarding path for all the flows in the system without exceeding the capacity of any of the edges, i.e., to completely avoid loss of packets from the given flows. **The source’s goal** is to progressively add flows, without exceeding the network’s capacity, forcing the controller to perform a flow swap. We shall show that the source has a strategy that forces the controller to swap traffic flows simultaneously in order to avoid packet loss.

Our model makes three basic assumptions: (i) each flow has a **fixed bandwidth**, (ii) the controller strives to **avoid packet loss**, and (iii) flows are **unsplittable**. We discuss these assumptions further in Sec. V.

The term *flow* in classic flow problems typically refers to the amount of traffic that is forwarded through each edge of the graph. Since our analysis focuses on SDN, we slightly divert from the common flow problem terminology, and use the term *flow* in its OpenFlow sense, i.e., a set of packets that share common properties, such as source and destination network addresses. A flow in our context, can be seen as a session between the source and destination that runs traffic at a fixed rate.

The network is represented by a directed weighted acyclic graph (Fig. 2b), $G = (\mathbb{V}, E, c)$, with a source s , a destination d , and a set of intermediate nodes, \mathbb{V}_{in} . Thus, $\mathbb{V} = \mathbb{V}_{in} \cup \{s, d\}$. The nodes directly connected to s are denoted by $\mathbb{O} = \{o_1, o_2, \dots, o_n\}$. Each of the outgoing edges from the source s has an infinite capacity, whereas the rest of the edges have a capacity c . For the sake of simplicity, and without loss of generality, throughout this section we assume that $c = 1$. Such a graph G is referred to as an *LFA graph*.

The source node progressively transmits traffic flows towards the destination node. Each flow represents a session between s and d ; every flow has a constant bandwidth, and cannot be split between two paths. A centralized controller configures the forwarding policy of the intermediate nodes, determining the path of each flow. Given a set of flows from s to d , the controller's goal is to configure the forwarding policy of the nodes in a way that allows all flows to be forwarded to d without exceeding the capacity of any of the edges.

The set of flows that are generated by s is denoted by $\mathbb{F} ::= \{F_1, F_2, \dots, F_k\}$. Each flow F_i is defined as $F_i ::= (i, f_i, r_i)$, where i is a unique flow index, f_i is the bandwidth satisfying $0 < f_i \leq c$, and r_i denotes the node that the controller forwards the flow to, i.e., $r_i \in \{o_1, o_2, \dots, o_n\}$.

It is assumed that the controller monitors the network, and thus it is aware of the flow set \mathbb{F} . The controller maintains a forwarding function, $R_{con} : \mathbb{F} \times \mathbb{V}_{in} \rightarrow \mathbb{V}_{in} \cup \{d\}$. Every node (switch) has a flow table, consisting of a set of *entries*; an element $w \in \mathbb{F} \times \mathbb{V}_{in}$ is referred to as an *entry* for short. An update of R_{con} is defined to be a partial function $u : \mathbb{F} \times \mathbb{V}_{in} \rightarrow \mathbb{V}_{in} \cup \{d\}$. We define a *reroute* as an update u that has a single entry in its domain. We call an update that has more than one entry in its domain a *swap*, and it is assumed that all updates in a *swap* are performed at the same time. We define a k -swap for $k \geq 2$ as a swap that updates entries in at least k different nodes. Note that a k -swap is possible only if $n \geq k$, where n is the number of nodes in \mathbb{O} . We focus our analysis on 2-swaps, and throughout the section we assume that $n \geq 2$. In Section II-F we discuss k -swaps for values of $k > 2$.

C. The LFA Game

The lossless flow allocation problem can be viewed as a game between two players, the source and the controller. The game proceeds by a sequence of steps; in each step the source

either adds or removes a single flow (Fig. 3), and then waits for the controller to perform a sequence of updates (Fig. 4). The source's strategy $\mathbb{S}_s(\mathbb{F}, R_{con}) = (a, F)$, is a function that defines for each flow set \mathbb{F} and forwarding function R_{con} for \mathbb{F} , a pair (a, F) representing the source's next step, where $a \in \{Add, Remove\}$ is the action to be taken by the source, and $F = (j, f_j, r_j)$ is a single flow to be added or removed. The controller's strategy is defined by $\mathbb{S}_{con}(R_{con}, a, F) = \mathbb{U}$, where $\mathbb{U} = \{u_1, \dots, u_\ell\}$ is a sequence of updates, such that (i) at the end of each update no edge exceeds its capacity, and (ii) at the end of the last update, u_ℓ , the forwarding function R_{con} defines a forwarding path for all flows in \mathbb{F} . Notice that when a flow is to be removed, the controller's update is trivial; it simply removes all the relevant entries from the domain of R_{con} . Hence our analysis focuses on *adding* new flows.

The following theorem, which is the crux of this section, argues that the source has a strategy that forces the controller to perform a swap, and thus that flow swaps are inevitable from the controller's perspective.

Theorem 1. *Let G be an LFA graph. In the LFA game over G , there exists a strategy, \mathbb{S}_s , for the source that forces every controller strategy, \mathbb{S}_{con} , to perform a 2-swap.*

Proof: Let m be the number of incoming edges to the destination node d in the LFA graph (see Fig 2b). For $m = 1$ the claim is trivial. Hence, we start by proving the claim for $m = 2$, i.e., there are two edges connected to node d , edges e_1 and e_2 . We show that the source has a strategy that, regardless of the controller's strategy, forces the controller to use a swap. In the first four steps of the game, the source generates four flows, $F_1 = (1, 0.35, o_1)$, $F_2 = (2, 0.35, o_1)$, $F_3 = (3, 0.45, o_2)$, and $F_4 = (4, 0.45, o_2)$, respectively. According to the Source Procedure of Fig. 3, after each flow is added, the source waits for the controller to update R_{con} before adding the next flow. After the flows are added, there are two possible cases:

- (a) The controller routes symmetrically through e_1 and e_2 , i.e. a flow of 0.35 and a flow of 0.45 through each of the edges. In this case the source's strategy at this point is to generate a new flow $F_5 = (5, 0.3, o_1)$ with a bandwidth of 0.3. The only way the controller can accommodate F_5 is by routing F_1 and F_2 through the same edge, allowing the new 0.3 flow to be forwarded through that edge. Since there is no sequence of *reroute* updates that allows the

SOURCE PROCEDURE

```

1   $\mathbb{F} \leftarrow \emptyset$ 
2  repeat at every step
3     $(a, F) \leftarrow \mathbb{S}_s(\mathbb{F}, R_{con})$ 
4    if  $a = Add$ 
5       $\mathbb{F} \leftarrow \mathbb{F} \cup F$ 
6      Wait for the controller to complete updates
7    else //  $a = Remove$ 
8       $\mathbb{F} \leftarrow \mathbb{F} \setminus F$ 
```

Fig. 3: The LFA game: the source's procedure.

CONTROLLER PROCEDURE

- 1 **repeat** at every step
 - 2 $\{u_1, \dots, u_\ell\} \leftarrow \mathbb{S}_{con}(R_{con}, a, F)$
 - 3 **for** $j \in [1, \ell]$
 - 4 Update R_{con} according to u_j
-

Fig. 4: The LFA game: the controller's procedure.

controller to reach the desired R_{con} , the only way to reach a state where F_1 and F_2 are routed through the same edge is to swap a 0.35 flow with a 0.45 flow. Thus, by issuing F_5 the controller forces a flow swap as claimed.

- (b) The controller routes F_1 and F_2 through one edge, and F_3 and F_4 through the other edge. In this case the source's strategy is to generate two flows, F_6 and F_7 , with a bandwidth of 0.2 each. The controller must route F_6 through the edge with F_1 and F_2 . Now each path sustains a bandwidth of 0.9 units. Thus, when F_7 is added by the source, the controller is forced to perform a swap between one of the 0.35 flows and one of the 0.45 flows.

In both cases the controller is forced to perform a 2-swap, swapping a flow from o_1 with a flow from o_2 . This proves the claim for $m = 2$.

The case of $m > 2$ is obtained by reduction to $m = 2$: the source first generates $m - 2$ flows with a bandwidth of 1 each, causing the controller to saturate $m - 2$ edges connected to node d (without loss of generality e_3, \dots, e_m). At this point there are only two available edges, e_1 and e_2 . From this point, the proof is identical to the case of $m = 2$. ■

The proof of Theorem 1 showed that the controller can be forced to perform a flow swap that involves $m = 2$ paths. For $m > 2$, we assumed that the source saturates $m - 2$ paths, reducing the analysis to the case of $m = 2$. In the following theorem we show that for $m > 2$ the controller can be forced to perform $\lfloor \frac{m}{2} \rfloor$ swaps.

Theorem 2. *Let G be an LFA graph. In the LFA game over G , if $m > 2$ then there exists a strategy, \mathbb{S}_s , for the source that forces every controller strategy, \mathbb{S}_{con} , to perform $\lfloor \frac{m}{2} \rfloor$ 2-swaps.*

Proof: Assume that m is even. The source generates m flows with a bandwidth of 0.35, m flows with a bandwidth of 0.45, and m flows with a bandwidth of 0.2. The only way the controller can route these flows without packet loss is as follows: each path sustains three flows with three different bandwidths, 0.2, 0.35, and 0.45. Now the source removes the m flows of 0.2, and adds $\frac{m}{2}$ flows of 0.3. As in case (a) of the proof of Theorem 1, adding each flow of 0.3 causes a 2-swap. The controller is thus forced to perform $\frac{m}{2} = \lfloor \frac{m}{2} \rfloor$ swaps.

If m is odd, then the source can saturate one of the edges by generating a flow with a bandwidth of 1, and then repeat the procedure above for the remaining $m - 1$ edges, yielding $\frac{m-1}{2} = \lfloor \frac{m}{2} \rfloor$ swaps. ■

For simplicity, throughout the rest of this section we assume that $m = 2$. However, as in Theorem 2, the analysis can be extended to the case of $m > 2$.

D. The Impact of Flow Swaps

We define a **metric** for flow swaps, by considering the oversubscription that is caused if the flows are **not** swapped simultaneously, but updated using an untimed approach.

We define the *oversubscription* of an edge, e , with respect to a forwarding function, R_{con} , to be the difference between the total bandwidth of the flows forwarded through e according to R_{con} , and the capacity of e . If the total bandwidth of the flows through e is less than the capacity of e , the oversubscription is defined to be zero.

Definition 1 (Flow swap impact). *Let \mathbb{F} be a flow set, and R_{con} be the corresponding forwarding function. Consider a 2-swap $u : \mathbb{F} \times \mathbb{V} \rightarrow \mathbb{V} \cup \{d\}$, such that $u = u_1 \cup u_2$, where $u_i = (w_i, v_i)$, for $w_i \in \mathbb{F} \times \mathbb{V}$, $v_i \in \mathbb{V} \cup \{d\}$, and $i \in \{1, 2\}$. The impact of u is defined to be the minimum of: (i) The oversubscription caused by applying u_1 to R_{con} , or (ii) the oversubscription caused by applying u_2 to R_{con} .*

Example 2. *We observe the scenario described in the proof of Theorem 1, and consider what would happen if the two flows had not been swapped simultaneously. The scenario had two cases; in the first case, the bandwidth through each edge is 0.8 before the controller swaps a 0.35 flow with a 0.45 flow. Thus, if the 0.35 flow is rerouted and then the 0.45 flow, the total bandwidth through the congested edge is $0.8 + 0.35 = 1.15$, creating a temporary oversubscription of 0.15. Thus, the flow swap impact in the first case is 0.15. In the second case, one edge sustains a bandwidth of 0.7, and the other a bandwidth of 0.9. The controller needs to swap a 0.35 flow with a 0.45 flow. If the controller first reroutes the 0.45 flow, then during the intermediate transition period, the congested edge sustains a bandwidth of $0.7 + 0.45 = 1.15$, and thus it is oversubscribed by 0.15. Hence, the impact in the second case is also 0.15.*

The following theorem shows that in the LFA game, the source can force the controller to perform a flow swap with a swap impact of roughly 0.5.

Theorem 3. *Let G be an LFA graph, and let $0 < \alpha < 0.5$. In the LFA game over G , there exists a strategy, \mathbb{S}_s , for the source that forces every controller strategy, \mathbb{S}_{con} , to perform a swap with an impact of α .*

Proof: Let $\epsilon = 0.1 - 0.2 \cdot \alpha$. We use the source's strategy from the proof of Theorem 1, with the exception that the bandwidths f_1, \dots, f_7 of flows F_1, \dots, F_7 are: $f_1 = f_2 = 0.5 - 2\epsilon$, $f_3 = f_4 = 0.5 - \epsilon$, $f_5 = 4\epsilon$, and $f_6 = f_7 = 3\epsilon$.

As in the proof of Theorem 1, there are two possible cases. In case (a), the controller routes symmetrically through the two paths, utilizing $1 - 3\epsilon$ of the bandwidth of each path. The source adds F_5 in response. To accommodate F_5 the controller swaps F_1 and F_3 . We determine the impact of this swap by considering the oversubscription of performing an untimed update; the controller first reroutes F_1 , and only then reroutes F_3 . Hence, the temporary oversubscription is $1 - 3\epsilon + 0.5 - 2\epsilon - 1 = 1.5 - 5\epsilon - 1$. Thus, the impact is $0.5 - 5\epsilon = \alpha$. In case (b), the controller forwards F_1 through the same path as F_2 , and F_3 through the same path as F_4 . The source responds by generating F_6 and F_7 . Again,

the controller is forced to swap between F_1 and F_3 . We compute the impact by considering an untimed update, where the controller reroutes F_3 first, causing an oversubscription of $1 - 4\epsilon + 0.5 - \epsilon - 1 = 0.5 - 5\epsilon = \alpha$. In both cases the source inflicts a flow swap with an impact of α . ■

Intuitively, Theorem 3 shows that not only are flow swaps inevitable, but they have a high impact on the network, as they can cause links to be congested by roughly 50% beyond their capacity.

E. Network Utilization

Theorem 1 demonstrates that regardless of the controller's policy, flow swaps cannot be prevented. However, the proof of Theorem 1 uses a scenario in which the edges leading to node d are almost fully utilized, suggesting that perhaps flow swaps are inevitable only when the traffic bandwidth is nearly equal to the max-flow of the graph. Arguably, as suggested in [23], by reserving some scratch capacity $\nu \cdot c$ through each of the edges, for $0 < \nu < 1$, it may be possible to avoid flow swaps. In the next theorem we show that if $\nu < \frac{1}{3}$, then flow swaps are inevitable.

Theorem 4. *Let G be an LFA graph, in which a scratch capacity of ν is reserved on each of the edges e_1, \dots, e_m , and let $\nu < \frac{1}{3}$. In the LFA game over G , there exists a strategy for the source, \mathbb{S}_s , that forces every controller strategy, \mathbb{S}_{con} , to perform a swap.*

Proof: We consider a graph G' , in which the capacity of each of the edges e_1, \dots, e_m is $1 - \nu$. By Theorem 3, for every $0 < \alpha < 0.5$, there exists a strategy for the source that forces a flow swap with an impact of α . Thus, there exists a strategy that forces at least one of the edges to sustain a bandwidth of $\alpha \cdot (1 - \nu)$. Since $\nu < \frac{1}{3}$, we have $(1 - \nu) > \frac{2}{3}$, and thus there exists an $\alpha < 0.5$ such that $\alpha \cdot (1 - \nu) > 1$. It follows that in the original graph G , with scratch capacity ν , there exists a strategy for the source that forces the controller to perform a flow swap in order to avoid the oversubscribed bandwidth of $\alpha \cdot (1 - \nu) > 1$. ■

The analysis of [23] showed that a scratch capacity of 10% is enough to address the reconfiguration scenarios that were considered in that work. Theorem 4 shows that even a scratch capacity of $33\frac{1}{3}\%$ does not suffice to prevent flow swaps scenarios. It follows that the 10% reserve that [23] suggest may not be sufficient in general for lossless reconfiguration.

F. n -Swaps

As defined above, a k -swap is a swap that involves k or more nodes. In previous subsections we discussed 2-swaps. The following theorem generalizes Theorem 1 to n -swaps, where n is the number of nodes in \mathbb{O} .

Theorem 5. *Let G be an LFA graph. In the LFA game over G , there exists a strategy, \mathbb{S}_s , for the source that forces every controller strategy, \mathbb{S}_{con} , to perform an n -swap.*

Proof: For $n = 1$, the claim is trivial. For $n = 2$, the claim was proven in Theorem 1. Thus, we assume $n \geq 3$.

If $m > 2$, the source first generates $m - 2$ flows with a rate c each, and we assume without loss of generality that after the controller allocates these flows only e_1 and e_2 remain unused. Thus, we focus on the case where $m = 2$.

We describe a strategy, \mathbb{S}_s as required; s generates three types of flows:

- Type A: two flows F_1, F_2 , at a rate of h each: $F_1 = (1, h, o_1)$, and $F_2 = (2, h, o_1)$.
- Type B: n flows, F_3, \dots, F_{n+2} , with a total rate g , i.e., at a rate of $\frac{g}{n}$ each. The source sends each of the n flows through a different node of \mathbb{O} .
- Type C: $n - 1$ flows, F_{n+3}, \dots, F_{2n+1} with a total rate g , i.e., $\frac{g}{n-1}$ each. The source sends each of the $n - 1$ flows through a different node of o_2, \dots, o_n .

We define h and g such that:

$$\frac{1}{3} < h < g < \frac{1}{2} \quad (1)$$

$$g > (n^2 - n)(1 - 2h) \quad (2)$$

We claim that for every n there exist g and h that satisfy (1) and (2). We prove this claim by finding g and h that satisfy the two conditions. We choose an arbitrary g in the range $(\frac{11}{24}, \frac{1}{2})$. We find a valid h by solving $g > (n^2 - n)(1 - 2h)$. The latter yields $h > \frac{1}{2} - \frac{g}{2(n^2 - n)}$. Since $n \geq 3$, we have $n^2 - n \geq 6$, and thus $\frac{g}{2(n^2 - n)} < \frac{0.5}{2 \times 6} = \frac{1}{24}$. Clearly, $\frac{g}{2(n^2 - n)} > 0$. It follows that every h that satisfies $\frac{1}{2} - \frac{1}{24} < h < \frac{1}{2} - 0$, also satisfies $h > \frac{1}{3}$. Hence, every g and h in the range $(\frac{11}{24}, \frac{1}{2})$ that satisfy $h < g$, also satisfy (1) and (2).

Intuitively, for h and g sufficiently close to $\frac{1}{2}$ (but less than $\frac{1}{2}$) (1) and (2) are satisfied.

We now prove that after generating the flows F_1, \dots, F_{2n+1} , the function R_{con} forwards all type B flows through the same path, and all type C flows through the same path. Assume by way of contradiction that there is a forwarding function R_{con} that forwards flows F_1, \dots, F_{2n+1} without loss, but does not comply to the latter claim. We consider two distinct cases: either the two type A flows are forwarded through the same edge, or they are forwarded through two different edges.

- If the two type A flows are forwarded through two different paths, then we assume that F_1 and the n type B flows are forwarded through e_1 and that F_2 and the $n - 1$ type C flows are forwarded through e_2 . Thus, at this point each of the two edges sustains traffic at a rate of $g + h$. By the assumption, there exists an update that swaps $i < n$ flows of type B with $j < n - 1$ flows of type C, such that after the swap none of the edges exceeds its capacity. Thus, the update adds the bandwidth $|j \cdot \frac{g}{n-1} - i \cdot \frac{g}{n}|$ to one of the edges, and this additional bandwidth must fit into the available bandwidth before the update, $1 - g - h$. Hence, $|j \cdot \frac{g}{n-1} - i \cdot \frac{g}{n}| < c - g - h$. Note that $1 - g - h < 1 - 2h < \frac{g}{n-1} - \frac{g}{n}$, following (1) and (2). Thus we get $|j \cdot \frac{g}{n-1} - i \cdot \frac{g}{n}| < \frac{g}{n-1} - \frac{g}{n}$. It follows that $|j \cdot n - i \cdot n + i| < 1$. Since j, i, n are integers, we get that $j \cdot n - i \cdot n + i = 0$, and thus $j = i \cdot \frac{n-1}{n}$. Now since $i \leq n$ and $j \leq n - 1$ are both natural numbers,

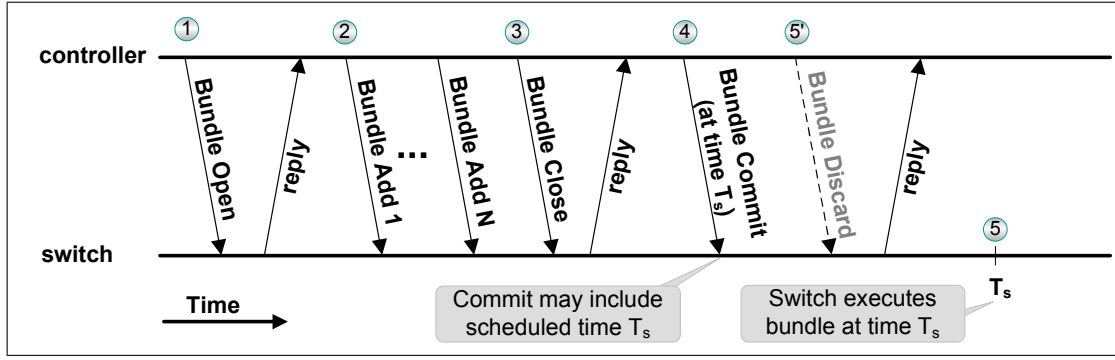


Fig. 5: A *Scheduled Bundle*: the *Bundle Commit* message may include T_s , the scheduled time of execution. The controller can use a *Bundle Discard* message to cancel the *Scheduled Bundle* before time T_s .

the only solution is $j = n - 1$ and $i = n$, which means that the flows from type B are all forwarded through the same path, as well as the flows of type C, contradicting the assumption.

- If the two type A flows are forwarded through the same edge, their total bandwidth is $2h$, and thus the remaining bandwidth through this edge is $1 - 2h$. From (2) we have $\frac{g}{n-1} - \frac{g}{n} > 1 - 2h$. We note that (i) $\frac{g}{n-1} > \frac{g}{n-1} - \frac{g}{n}$, and (ii) $\frac{g}{n} > \frac{g}{n-1} - \frac{g}{n}$. It follows that $\frac{g}{n-1} > 1 - 2h$, and also $\frac{g}{n} > 1 - 2h$, and thus none of the type B or type C flows fit on the same path with F_1 and F_2 . Thus, all the type B and type C flows are on the same path, contradicting the assumption.

We have shown that all flows of type B, denoted by \mathbb{F}^B , must be forwarded through the same path, and that all flows of type C, denoted by \mathbb{F}^C , are forwarded through the same path. Thus, after the source generates the $2 \cdot n + 1$ flows, there are two possible scenarios:

- The two type A flows are forwarded through the same path, and the type B and type C flows are forwarded through the other path. In this case s generates two flows at a rate of $1 - h - g$ each. To accommodate both flows the controller must swap the flows of \mathbb{F}^B with F_1 or the flows of \mathbb{F}^C with F_2 . Both possible swaps involve n entries, and thus the controller is forced to perform an n -swap.
- One path is used for F_1 and the flows of \mathbb{F}^C , and the other path is used for F_2 and the flows of \mathbb{F}^B . In this case the source generates a flow with a bandwidth of $1 - 2h$, again forcing the controller to swap the flows of \mathbb{F}^B with F_1 or the flows of \mathbb{F}^C with F_2 .

In both cases the controller is forced to perform a swap that involves the n nodes, i.e., an n -swap. ■

III. DESIGN AND IMPLEMENTATION

A. Protocol Design

1) Overview

A TIME4-enabled system is comprised of two main components:

- **OpenFlow time extension.** TIME4 is built upon the OpenFlow protocol. We define an extension to the OpenFlow protocol that enables timed updates; the controller

can attach an *execution time* to every OpenFlow command it sends to a switch, defining when the switch should perform the required command. It should be noted that the TIME4 approach is not limited to OpenFlow; we have defined a similar time extension to the NETCONF protocol [45], but in this paper we focus on TIME4 in the context of OpenFlow, as described in the next subsection.

- **Clock synchronization.** TIME4 requires the switches and controller to maintain a local clock, allowing time-triggered events. Hence, the local clocks should be synchronized. The OpenFlow time extension we defined does not mandate a specific synchronization method. Various mechanisms may be used, e.g., the Network Time Protocol (NTP), the Precision Time Protocol (PTP) [5], or GPS-based synchronization. The prototype we designed and implemented uses REVERSEPTP [46], as described below.

2) OpenFlow Time Extension

We present an extension that allows OpenFlow controllers to signal the time of execution of a command to the switches. This extension is described in full in Appendix A.²

Our extension makes use of the OpenFlow [19] **Bundle** feature; a Bundle is a sequence of OpenFlow messages from the controller that is applied as a single operation. Our time extension defines *Scheduled Bundles*, allowing all commands of a Bundle to come into effect at a pre-determined time. This is a generic means to extend all OpenFlow commands with the scheduling feature.

Using Bundle messages for implementing TIME4 has two significant advantages: (i) It is a generic method to add the time extension to all OpenFlow commands without changing the format of all OpenFlow messages; only the format of Bundle messages is modified relative to the Bundle message format in [19], optionally incorporating an execution time. (ii) The Scheduled Bundle allows a relatively straightforward way to *cancel* scheduled commands, as described below.

Fig. 5 illustrates the *Scheduled Bundle* message procedure. In step 1, the controller sends a *Bundle Open* message to the switch, followed by one or more Add messages (step 2). Every *Add* message encapsulates an OpenFlow message, e.g.,

²A preliminary version of this extension was presented in [47].

a *FLOW_MOD* message. A *Bundle Close* is sent in step 3, followed by the *Bundle Commit* (step 4), which optionally includes the scheduled time of execution, T_s . The switch then executes the desired command(s) at time T_s .

The *Bundle Discard* message (step 5') allows the controller to enforce an all-or-none scheduled update; after the *Bundle Commit* is sent, if one of the switches sends an *error* message, indicating that it is unable to schedule the current Bundle, the controller can send a Discard message to all switches, canceling the scheduled operation. Hence, when a switch receives a scheduled commit, to be executed at time T_s , the switch can verify that it can dedicate the required resources to execute the command as close as possible to T_s . If the switch's resources are not available, for example due to another command that is scheduled to T_s , then the switch replies with an error message, aborting the scheduled commit. Significantly, this mechanism allows switches to execute the command with a guaranteed scheduling accuracy, avoiding the high variation that occurs when untimed updates are used.

The OpenFlow time extension also defines *Bundle Feature Request* messages, which allow the controller to query switches about whether they support Scheduled Bundles, and to configure some of the switch parameters related to *Scheduled Bundles*.

3) Clock Synchronization: REVERSEPTP

In the last decade PTP, based on the IEEE 1588 [5] standard, has become a common feature in commodity switches, typically providing a clock accuracy on the order of 1 microsecond.

In [46], [48] we introduced REVERSEPTP a PTP variant for SDNs. REVERSEPTP is based on PTP, but is conceptually reversed. In PTP a single node periodically distributes its time to the other nodes in the network. In REVERSEPTP all nodes in the network (the switches) periodically distribute their time to a single node (the controller). The controller keeps track of the offsets, denoted by offset_i for switch i , between its clock and each of the switches' clocks, and uses them to send each switch individualized timed commands.

REVERSEPTP allows the complex clock algorithms to be implemented by the controller, whereas the 'dumb' switches only need to distribute their time to the controller. Following the SDN paradigm, the REVERSEPTP algorithmic logic can be programmed and dynamically tuned at the controller without affecting the switches.

Another advantage of REVERSEPTP, which played an important role in our experiments, is that REVERSEPTP allows the controller to keep track of the synchronization status of each clock; a clock synchronization protocol requires a long setup time, typically tens of minutes. REVERSEPTP provides an indication of when the setup process has completed.

As shown in [46], REVERSEPTP can be effectively used to perform timed updates; in order to have switch i perform a command at time T_s , the controller instructs i to perform the command at time T_s^i , where $T_s^i = T_s + \text{offset}_i$ takes the offset

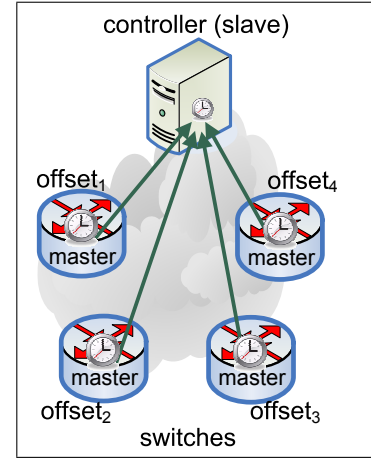


Fig. 6: REVERSEPTP in SDN: switches distribute their time to the controller. Switches' clocks are *not* synchronized. For every switch i , the controller knows offset_i between switch i 's clock and its local clock.

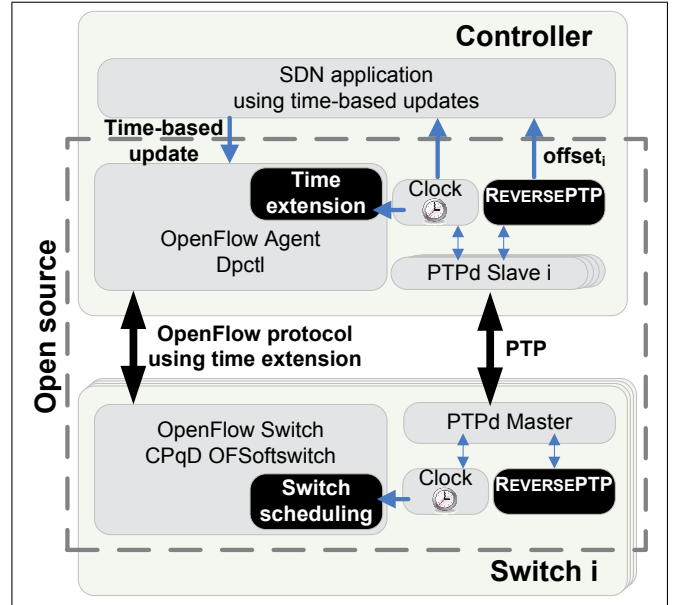


Fig. 7: TIME4 prototype design: the black blocks are the components implemented in the context of this work.

between the controller and switch i into account,³ causing i to perform the action at time T_s according to the controller's clock.

B. Prototype Design and Implementation

We have designed and implemented a software-based prototype of TIME4, as illustrated in Fig. 7. The components we implemented are marked in black. These components run on Linux, and are publicly available as open source [43].

Our TIME4-enabled OFSoftswitch prototype was adopted

³ T_s^i , as described above is a first order approximation of the desired execution time. The controller can compute a more accurate execution time by also considering the clock skew and drift, as discussed in [46].

by the ONF as the official prototype of Scheduled Bundles.⁴

Switches. Every switch i runs an OpenFlow switch software module. Our prototype is based on the open source CPqD OF-Softswitch [49],⁵ incorporating the *switch scheduling* module (see Fig. 7) that we implemented. When the switch receives a *Scheduled Bundle* from the controller, the *switch scheduling* module schedules the respective OpenFlow command to the desired time of execution. The switch scheduling module also handles *Bundle Feature Request* messages received from the controller.

Each switch runs a REVERSEPTP master, which distributes the switch's time to the controller. Our REVERSEPTP prototype is a lightweight set of Bash scripts that is used as an abstraction layer over the well-known open source PTPd [50] module. Our software-based implementation uses the Linux clock as the reference for PTPd, and for the switch's scheduling module. To the best of our knowledge, ours is the first open source implementation of REVERSEPTP.

Controller. The controller runs an OpenFlow agent, which communicates with the switches using the OpenFlow protocol. Our prototype uses the CPqD Dpctl (Datapath Controller), which is a simple command line tool for sending OpenFlow messages to switches. We have extended Dpctl by adding the time extension; the Dpctl command-line interface allows the user to define the execution time of a *Bundle Commit*. Dpctl also allows a user to send a *Bundle Feature Request* to switches.

The controller runs REVERSEPTP with n instances of PTPd in slave mode, where n is the number of switches in the network. One or more SDN applications can run on the controller and perform timed updates. The application can extract the offset, offset_i , of every switch i from REVERSEPTP, and use it to compute the scheduled execution time of switch i in every timed update. The Linux clock is used as a reference for PTPd, and for the SDN application(s).

IV. EVALUATION

A. Evaluation Method

Environment. We evaluated our prototype on a 71-node testbed in the DeterLab [51] environment. Each machine (PC) in the testbed either played the role of an OpenFlow switch, running our TIME4-enabled prototype, or the role of a host, sending and receiving traffic. A separate machine was used as a controller, which was connected to the switches using an out-of-band network.

We remark that we did not use Mininet [52] in our evaluation, as Mininet is an emulation environment that runs on a single machine, making it impractical for emulating simultaneous or time-triggered events. We did, however, run our prototype over Mininet in some of our preliminary testing and verification.

⁴The ONF process for adding new features to OpenFlow requires every new feature to be prototyped.

⁵OFSoftswitch is one of the two software switches used by the Open Networking Foundation (ONF) for prototyping new OpenFlow features. We chose this switch since it was the first open source OpenFlow switch to include the Bundle feature.

Performance attributes. Three performance attributes play a key role in our evaluation, as shown in Table I.

Δ	The average time elapsed between two consecutive messages sent by the controller.
I_R	Installation latency range: the difference between the maximal rule installation latency and the minimal installation latency.
δ	Scheduling error: the maximal difference between the actual update time and the scheduled update time.

TABLE I: Performance Attributes.

Intuitively, Δ and I_R determine the performance of untimed updates. Δ indicates the **controller's** performance; an OpenFlow controller can handle as many as tens of thousands [53] to millions [54] of packets per second, depending on the type of controller and the machine's processing power. Hence, Δ can vary from 1 microsecond to several milliseconds. I_R indicates the installation latency variation. The installation latency is the time elapsed from the instant the controller sends a rule modification message until the rule has been installed. The installation latency of an OpenFlow rule modification (FLOW_MOD) has been shown to range from 1 millisecond to seconds [28], [55], and grows dramatically with the number of installations per second.

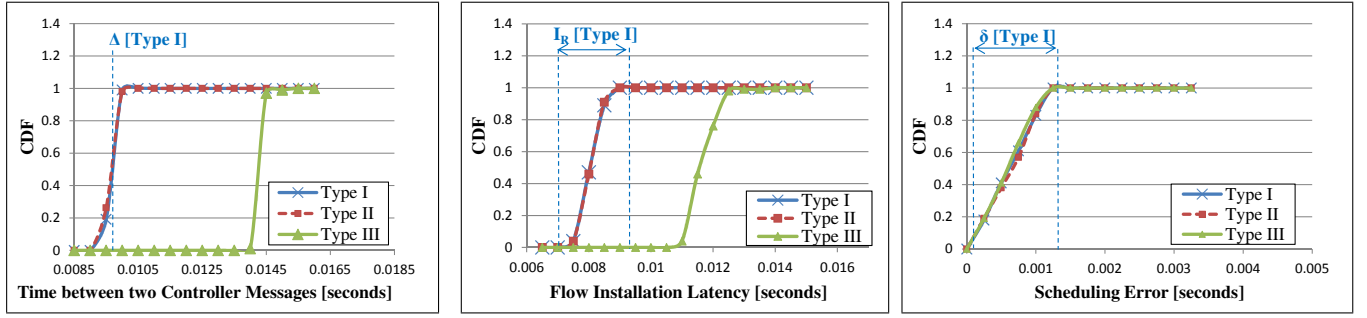
The attribute that affects the performance of **timed** updates is the switches' scheduling error, δ . When an update is scheduled to be performed at time T_0 , it is performed in practice at some time $t \in [T_0, T_0 + \delta]$.⁶ The scheduling error, δ , is affected by two factors: the device's *clock accuracy*, which is the maximal offset between the clock value and the value of an accurate time reference, and the *execution accuracy*, which is a measure of how accurately the device can perform a timed update, given run-time parameters such as the concurrently executing tasks and the load on the device. The achievable clock accuracy strongly depends on the network size and topology, and on the clock synchronization method. For example, the clock accuracy using the Precision Time Protocol [5] is typically on the order of 1 microsecond (e.g., [6]).

Software-based evaluation. Our experiments measure the three performance attributes in a setting that uses **software switches**. While the *values* we measured do not necessarily reflect on the performance of systems that use hardware-based switches, the merit of our evaluation is that we **vary** these parameters and analyze how they affect the network update performance with untimed approaches and with TIME4.

B. Performance Attribute Measurement

Our experiments measured the three attributes, Δ , I_R , and δ , illustrating how accurately updates can be applied in software-based OpenFlow implementations. It should be noted that these three values depend on the processing power of the testbed machine; we measured the parameters for three types of DeterLab machines, Type I, II, and III, listed in Table II. Each attribute was measured 100 times on each machine type, and Fig. 8 illustrates our results. The figure graphically depicts the values Δ , I_R , and δ of machine Type I as an example.

⁶An alternative representation of the accuracy, δ , assumes a symmetric error, $T_0 \pm \delta$. The two approaches are equivalent.



(a) The empirical Cumulative Distribution Function (CDF) of the time elapsed between two consecutive controller messages. Δ is the average value, which is shown in the figure for Type I. (b) The empirical CDF of the flow installation latency. I_R is the difference between the max and min values, as shown in the figure for Type I. (c) The empirical CDF of the scheduling error, i.e., the difference between the actual execution time and the scheduled execution time. δ is the maximal error value, as shown in the figure for Type I.

Fig. 8: Measurement of the three performance attributes: (a) Δ , (b) I_R , and (c) δ .

Machine Type		Δ	I_R	δ
I	Intel Xeon E3 LP 2.4 GHz, 16 GB RAM	9.64	1.3	1.23
II	Intel Xeon 2.1 GHz, 4 GB RAM	9.6	1.47	1.18
II	Intel Dual Xeon 3 GHz, 2 GB RAM	14.27	2.72	1.19

TABLE II: Measured attributes in milliseconds.

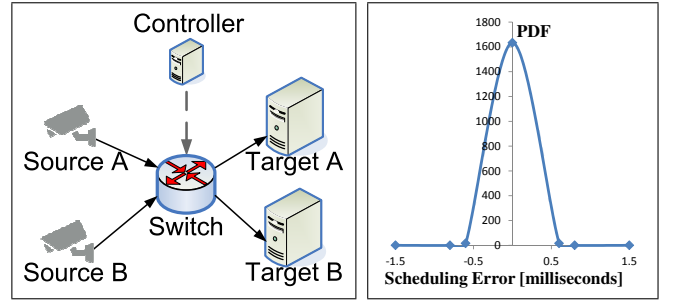
The measured scheduling error, δ , was slightly more than 1 millisecond in all the machines we tested. Our experiments showed that the *clock accuracy* using REVERSEPTP over the DeterLab testbed is on the order of 100 microseconds. The measured value of δ in Table II shows the *execution accuracy*, which is an order of magnitude higher. The installation latency range, I_R , was slightly higher than δ , around 1 to 3 milliseconds. The measured value of Δ was high, on the order of 10 milliseconds, as Dpctl is not optimized for performance.

In software-based switches, the CPU handles both the data-plane traffic and the communication with the controller, and thus I_R and δ can be affected by the rate of data-plane traffic through the switch. Hence, in our experiments we fixed the rate of traffic through each switch to 10 Mbps, allowing an ‘apples-to-apples’ comparison between experiments.

C. Microbenchmark: Video Swapping

To demonstrate how TIME4 is used in a real-life scenario, we reconstructed the video swapping topology of [29], as illustrated in Fig. 9a. Two video cameras, A and B, transmit an uncompressed video stream to targets A and B, respectively. At a given point in time, the two video streams are swapped, so that the stream from source A is transmitted to target B, and the stream from B is sent to target A. As described in [29], the swap must be performed at a specific time instant, in which the video sources transmit data that is not visible to the viewer, making the swap unnoticeable.

The authors of [29] noted that the precisely-timed swap cannot be performed by an OpenFlow switch, as currently OpenFlow does not provide abstractions for performing accurately timed changes. Instead, it uses *source timing*, where



(a) Topology. (b) Video swapping accuracy.

Fig. 9: Microbenchmark: video swapping.

sources A and B are time-synchronized, and determine the swap time by using a swap indication in the packet header. The OpenFlow switch acts upon the swap indication to determine the correct path for each stream. We note that the main drawback of this source-timed approach is that the SMPTE 2022-6 video streaming standard [56], which was used in [29], does not currently define an indication about where in the video stream a packet comes from, and specifically does not include an indication about the correct swapping time. Hence, off-the-shelf streaming equipment does not provide this indication. In [29], the authors used a dedicated Linux server to integrate the non-standard swap indication.

In this experiment we studied how TIME4 can tackle the video swapping scenario, avoiding the above drawback. Each node in the topology of Fig. 9a was emulated by a DeterLab machine. We used two 10 Mbps flows, generated by Iperf [57], to simulate the video streams. Each swap was initiated by the controller 100 milliseconds in advance (as in [29]): the controller sent a Scheduled Bundle, incorporating two updates, one for each of the flows. We repeated the experiment 100 times, and measured the scheduling error.

The measurement was performed by analyzing capture files taken at the sources and at the switch’s egress ports. A swap that was scheduled to be performed at time T , was considered accurate if every packet that was transmitted by each of the source before time T was forwarded according to the old

configuration, and every packet that was transmitted after T was forwarded according to the new configuration. The scheduling error of each swap (measured in milliseconds) was computed as the number of misrouted packets, divided by the bandwidth of the traffic flow. The sign of the scheduling error indicates whether the swap was performed before the scheduled time (negative error) or after it (positive error).

Fig. 9b illustrates the empirical Probability Density Function (PDF) of the scheduling error of the swap, i.e., the difference between the actual swapping time and the scheduled swapping time. As shown in the figure, the swap is performed within ± 0.6 milliseconds of the scheduled swap time. We note that this is the achievable accuracy in a software-based OpenFlow switch, and that a much higher degree of accuracy, on the order of microseconds, can be achieved if two conditions are met: (i) A hardware switch is used, supporting timed updates with a microsecond accuracy, as shown in [30], and (ii) The cameras are connected to the switch over a single hop, allowing low latency variation, on the order of microseconds.

D. Flow Swap Evaluation

1) Experiment Setting

We evaluated our prototype on a 71-node testbed under. We used the testbed to emulate an OpenFlow network with 32 hosts and 32 leaf switches, as depicted in Fig. 11, with $n = 32$.

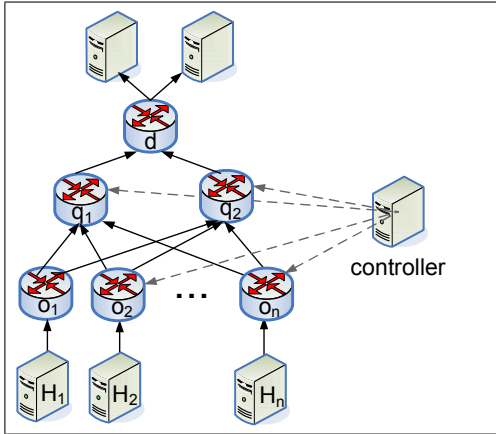


Fig. 11: Experimental evaluation: every host and switch was emulated by a Linux machine in the DeterLab testbed. All links have a capacity of 10 Mbps. The controller is connected to the switches by an out-of-band network.

Metric. A flow swap that is not performed in a coordinated way may bare a high cost: either packet loss, deep buffering, or a combination of the two. We use packet loss as a **metric** for the cost of flow swaps, assuming that deep buffering is not used.

We used Iperf to generate flows from the sources to the destination, and to measure the number of packets lost between the source and the destination.

The flow swap scenario. All experiments were flow swaps with a swap impact of 0.5.⁷ We used two static flows,

which were not reconfigured in the experiment: H_1 generates a 5 Mbps flow that is forwarded through q_1 , and H_2 generates a 5 Mbps flow that is generated through q_2 . We generated n additional flows (where n is the number of switches at the bottom layer of the graph): (i) A 5 Mbps flow from H_1 to the destination. (ii) $n - 1$ flows, each having a bandwidth of $\frac{5}{n-1}$ Mbps. Every flow swap in our experiment required the flow of (i) to be swapped with the $n - 1$ flows of (ii). Note that this swap has an impact of 0.5.

2) Experimental Results

TIME4 vs. other update approaches. In this experiment we compared the packet loss of TIME4 to other update approaches described in Sec. I-D. As discussed in Sec. I-D, applying the *order* approach or the *two-phase* approach to flow swaps produces similar results. This observation is illustrated in Fig. 10b. In the rest of this section we refer to these two approaches collectively as the *untimed* approaches.

In our experiments we also implemented a SWAN-based [23] update, and a B4-based [24] update. In SWAN, we used a 10% scratch on each of the links, and in B4 updates we temporarily reduced the bandwidth of each flow by 10% to avoid packet loss. As depicted in Fig. 10b, SWAN and B4 yield a slightly lower packet loss rate than TIME4; the average number of packets lost in each TIME4 flow swap is 0.2, while with SWAN and B4 only 0.1 packets are lost on average.

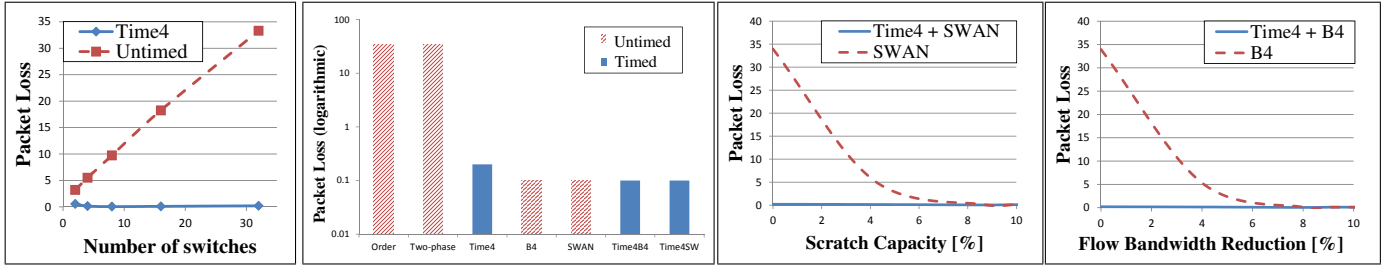
To study the effect of using **time** in SWAN and in B4, we also performed hybrid updates, illustrated in Fig. 10c and 10d, and in the two right-most bars of Fig. 10b. We combined SWAN and TIME4, by performing a timed update on a network with scratch capacity, and compared the packet loss to the conventional SWAN-based update. We repeated the experiment for various values of scratch capacity, from 0% to 10%. As illustrated in Fig. 10c, the TIME4+SWAN approach can achieve the same level of packet loss as SWAN with **less scratch capacity**. We performed a similar experiment with a timed B4 update, varying the bandwidth reduction rate between 0% and 10%, and observed similar results.

Number of switches. We evaluated the effect of n , the number of switches involved in the flow swap, on the packet loss. We performed an n -swap with $n = 2, 4, 8, 16, 32$. As illustrated in Fig. 10a, the number of packets lost during an untimed update grows linearly with the number of switches n , while the number of packets lost in a TIME4 update is less than one on average, and is not affected by the number of switches. As n increases, the update duration⁸ is longer, and hence more packets are lost during the update procedure.

Controller performance. In this experiment we explored how the controller's performance, represented by Δ , affects the packet loss rate in an untimed update. As Δ increases, the update procedure requires a longer period of time, and hence more packets are lost (Fig. 12) during the process. We note that although previous work has shown that Δ can be on the order of microseconds in some cases [54], Dpctl is not optimized for performance, and hence Δ in our experiments was on the

⁷By Theorem 3, the source can force the controller to perform a flow swap with an impact as high as roughly 0.5.

⁸The **update duration** is the time elapsed from the instant the first switch is updated until the instant the last switch is updated. In our setting the update duration is roughly $(n - 1)\Delta$.



(a) The no. of packets lost in a flow swap vs. no. of switches involved in the update. (b) The number of packets lost in a flow swap in different update approaches (with $n = 32$). (c) The number of packets lost in a flow swap using SWAN and TIME4 + SWAN (with $n = 32$). (d) The number of packets lost in a flow swap using B4 and TIME4 + B4 (with $n = 32$).

Fig. 10: Flow swap performance: in large networks (a) TIME4 allows significantly less packet loss than untimed approaches. The packet loss of TIME4 is slightly higher than SWAN and B4 (b), while the latter two methods incur higher overhead. Combining TIME4 with SWAN or B4 provides the best of both worlds; low packet loss (b) and low overhead (c and d).

order of milliseconds. As shown in Fig. 12, we synthetically increased Δ , and observed its effect on the packet loss during flow swaps.

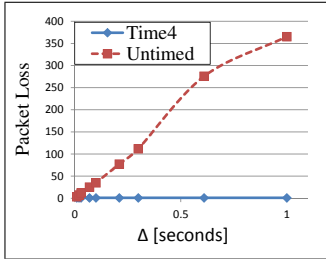
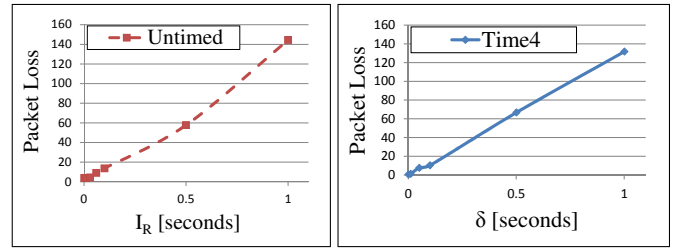


Fig. 12: The number of packets lost in a flow swap vs. Δ . The packet loss in TIME4 is not affected by the controller's performance (Δ).

Installation latency variation. Our next experiment (Fig. 13a) examined how the installation latency variation, denoted by I_R , affects the packet loss during an untimed update. We analyzed different values of I_R : in each update we synthetically determined a uniformly distributed installation latency, $I \sim U[0, I_R]$. As shown in Fig. 13a, the switch's installation latency range, I_R , dramatically affects the packet loss rate during an untimed update. Notably, when I_R is on the order of 1 second, as in the extreme scenarios of [28], [55], TIME4 has a significant advantage over the untimed approach.

Scheduling error. Figure 13b depicts the packet loss as a function of the scheduling error of TIME4. By Fig. 10a, 13a and 13b, we observe that if δ is sufficiently low compared to I_R and $(n - 1)\Delta$, then TIME4 outperforms the untimed approaches. Note that even if switches are not implemented with extremely low scheduling error δ , we expect TIME4 to outperform the untimed approach, as typically $\delta < I_R$, as further discussed in Section V.

Summary. The experiments presented in this section demonstrate that TIME4 performs significantly better than untimed approaches, especially when the update involves multiple switches, or when there is a non-deterministic installation latency. Interestingly, TIME4 can be used in conjunction with existing approaches, such as SWAN and B4, allowing the same level of packet loss **with less overhead** than the untimed



(a) The number of packets lost in a flow swap vs. the installation latency range, I_R . (b) The number of packets lost in a flow swap vs. the scheduling error, δ .

Fig. 13: Performance as a function of I_R and δ . Untimed updates are affected by the installation latency variation (I_R), whereas TIME4 is affected by the scheduling error (δ). TIME4 is advantageous since typically $\delta < I_R$.

variants.

V. DISCUSSION

1) Scheduling accuracy

The advantage of timed updates greatly depends on the **scheduling accuracy**, i.e., on the switches' ability to accurately perform an update at its scheduled time. Clocks can typically be synchronized on the order of 1 microsecond (e.g., [6]) using PTP [5]. However, a switch's ability to accurately perform a scheduled action depends on its implementation.

- *Software switches:* Our experimental evaluation showed that the scheduling error in the software switches we tested was on the order of 1 millisecond.
- *Hardware-based scheduling:* The work of [30] has shown a method that allows the scheduling error of timed events in hardware switches to be as low as 1 microsecond.
- *Software-based scheduling in hardware switches:* A scheduling mechanism that relies on the switch's software may be affected by the switch's operating system and by other running tasks. Measures can be taken to implement an accurate software-based scheduling in TIME4: when a switch is aware of an update that is scheduled to take place at time T_s , it can avoid performing heavy maintenance tasks at this time, such as TCAM entry

rearrangement. Update messages received slightly before time T_s can be queued and processed after the scheduled update is executed. Moreover, if a switch receives a timed command that is scheduled to take place at the same time as a previously received command, it can send an error message to the controller, indicating that the last received command cannot be executed.

It is an important observation that in a typical system we expect the scheduling error to be lower than the installation latency variation, i.e., $\delta < I_R$. Untimed updates have a non-deterministic installation latency. On the other hand, timed updates are predictable, and can be scheduled in a way that avoids conflicts between multiple updates, allowing δ to be typically lower than I_R .

2) Model assumptions

Our model assumes a *lossless* network with *unsplittable*, *fixed-bandwidth* flows. A notable example of a setting in which these assumptions are often valid is a WAN or a carrier network. In carrier networks the maximal **bandwidth** of a service is defined by its bandwidth profile [27]. Thus, the controller cannot dynamically change the bandwidth of the flows, as they are determined by the SLA. The Frame Loss Ratio (FLR) is one of the key performance attributes [27] that a service provider must comply to, and cannot be compromised. **Splitting** a flow between two or more paths may result in packets being received out-of-order. Packet reordering is a key performance parameter in carrier-grade performance and availability measurement, as it affects various applications such as real-time media streaming [58]. Thus, all packets of a flow are forwarded through the same path.

3) Short term vs. long term scheduling

The OpenFlow time extension we presented in Section III is intended for short term scheduling; a controller should schedule an action to a near-future time, on the order of seconds in the future. The challenge in long term scheduling is that during the long period between the time at which the Scheduled Bundle was sent and the time at which it is meant to be executed various external events may occur: the controller may fail or reboot, or a second controller⁹ may try to perform a conflicting update. Near future scheduling guarantees that external events that may affect the scheduled operation such as a switch reboot have a low probability of occurring. Since near-future scheduling is on the order of seconds, this short potentially hazardous period is no worse than in conventional updates, where an OpenFlow command may be executed a few seconds after it was sent by the controller.

4) Network latency

In Fig. 1, the switches S_1 and S_3 are updated at the same time, as it is implicitly assumed that all the links have the same latency. In the general case each link has a different latency, and thus S_1 and S_3 should not be updated at the same time, but at two different times, T_1 and T_3 , that account for the different latencies.

⁹In an SDN with a distributed control plane, where more than one controller is used.

5) Failures

A timed update may fail to be performed in a coordinated way at multiple switches if some of the switches have failed, or if some of the controller commands have failed to reach some of the switches. Therefore, the controller uses a reliable transport protocol (TCP), in which dropped packets are re-transmitted. If the controller detects that a switch has failed, or failed to receive some of the Bundle messages, the controller can use the *Bundle Discard* to cancel the coordinated update. Note that the controller should send timed update messages sufficiently ahead of the scheduled time of execution, allowing enough time for possible retransmission and Discard message transmission.

6) Controller performance overhead

The prototype design we presented (Fig. 7) uses REVERSEPTP [46] to synchronize the switch and the controllers. A synchronization protocol may yield some performance overhead on the controller and switches, and some overhead on the network bandwidth. In our experiments we observed that the CPU utilization of the PTP processes in the controller in an experiment with 32 switches was 5% on the weakest machine we tested, and significantly less than 1% on the stronger machines. As for the network bandwidth overhead, accurate synchronization using PTP typically requires the controller to exchange ~ 5 packets per second per switch [59], a negligible overhead in high-speed networks.

VI. CONCLUSION

Time and clocks are valuable tools for coordinating updates in a network. We have shown that dynamic traffic steering by SDN controllers requires flow swaps, which are best performed as close to instantaneously as possible. Time-based operation can help to achieve carrier-grade packet loss rate in environments that require rapid path reconfiguration. Our OpenFlow time extension can be used for implementing flow swaps and TIME4. It can also be used for a variety of additional timed update scenarios that can help improve network performance during path and policy updates.

VII. ACKNOWLEDGMENTS

We gratefully acknowledge Oron Anschel and Nadav Shiloach, who implemented the TIME4-enabled OFSoftswitch prototype. We thank Jean Tourrilhes and the members of the Extensibility working group of the ONF for many helpful comments that contributed to the OpenFlow time extension. We also thank Nate Foster, Laurent Vanbever, Joshua Reich and Isaac Keslassy for helpful discussions. We gratefully acknowledge the DeterLab project [51] for the opportunity to perform our experiments on the DeterLab testbed. This work was supported in part by the ISF grant 1520/11.

REFERENCES

- [1] T. Mizrahi and Y. Moses, "Software defined networks: It's about time," in *IEEE INFOCOM*, 2016.
- [2] T. Mizrahi and Y. Moses, "TIME4: Time for SDN," technical report, arXiv preprint arXiv:1505.03421, 2015.
- [3] ITU-T G.8271/Y.1366, "Time and phase synchronization aspects of packet networks," *ITU-T*, 2012.
- [4] J. C. Corbett, J. Dean, M. Epstein, A. Fikes, C. Frost, J. Furman, S. Ghemawat, A. Gubarev, C. Heiser, P. Hochschild, et al., "Spanner: Google's globally distributed database," *ACM Transactions on Computer Systems (TOCS)*, vol. 31, no. 3, p. 8, 2013.
- [5] IEEE TC 9, "1588 IEEE Standard for a Precision Clock Synchronization Protocol for Networked Measurement and Control Systems Version 2," *IEEE*, 2008.
- [6] H. Li, "IEEE 1588 time synchronization deployment for mobile backhaul in China Mobile," keynote presentation, IEEE ISPCS, 2014.
- [7] IEEE Std C37.238, "IEEE Standard Profile for Use of IEEE 1588 Precision Time Protocol in Power System Applications," *IEEE*, 2011.
- [8] "ONF SDN Product Directory," <https://www.opennetworking.org/sdn-resources/onf-products-listing>, January, 2015.
- [9] "Broadcom BCM56840+ Switching Technology," product brief, http://www.broadcom.com/collateral/pb/56840_PLUS-PB00-R.pdf, 2011.
- [10] "Broadcom BCM56850 StrataXGS Trident II Switching Technology," product brief, <http://www.broadcom.com/collateral/pb/56850-PB03-R.pdf>, 2013.
- [11] "Centec CTC6048 Advanced Product Brief," http://www.valleytalk.org/wp-content/uploads/2010/12/CTC6048-Product-Brief_v2.5.pdf, 2010.
- [12] "Intel Ethernet Switch FM5000/FM6000 Datasheet," <http://www.intel.com/content/dam/www/public/us/en/documents/datasheets/ethernet-switch-fm5000-fm6000-datasheet.pdf>, 2014.
- [13] "LSI/Avago Axxia Communication Processor AXM5500 Family," https://www.lsi.com/downloads/Public/Communication%20Processors/Axxia%20Communication%20Processor/LSI_PB_AXM5500_E.pdf, 2014.
- [14] "Mellanox SwitchX-2 Product Brief," http://www.mellanox.com/related-docs/prod_silicon/SwitchX-2_EN_SDN.pdf, 2013.
- [15] "Tilera TILE-Gx8072 Processor Product Brief," http://www.tilera.com/sites/default/files/productbriefs/TILE-Gx8072_PB041-04_WEB.pdf, 2014.
- [16] "Marvell ARMADA XP Functional Spec," <http://www.marvell.com/embedded-processors/armada-xp/assets/ARMADA-XP-Functional-SpecDatasheet.pdf>, 2014.
- [17] "Marvell Xelerated HX4100," product brief, http://www.marvell.com/network-processors/assets/Xelerated_HX4100-02_product%20brief_v8.pdf, 2014.
- [18] N. McKeown, T. Anderson, H. Balakrishnan, G. Parulkar, L. Peterson, J. Rexford, S. Shenker, and J. Turner, "Openflow: enabling innovation in campus networks," *SIGCOMM Comput. Commun. Rev.*, vol. 38, pp. 69–74, Mar. 2008.
- [19] Open Networking Foundation, "Openflow switch specification," *Version 1.4.0*, 2013.
- [20] "Interface to the Routing System (I2RS) working group," <https://datatracker.ietf.org/wg/i2rs/charter/>, IETF, 2016.
- [21] "Forwarding and Control Element Separation (ForCES) working group," <https://datatracker.ietf.org/wg/forces/charter/>, IETF, 2016.
- [22] M. Al-Fares, S. Radhakrishnan, B. Raghavan, N. Huang, and A. Vahdat, "Hedera: Dynamic flow scheduling for data center networks," in *NSDI*, 2010.
- [23] C.-Y. Hong, S. Kandula, R. Mahajan, M. Zhang, V. Gill, M. Nanduri, and R. Wattenhofer, "Achieving high utilization with software-driven WAN," in *ACM SIGCOMM*, 2013.
- [24] S. Jain et al., "B4: Experience with a globally-deployed software defined WAN," in *ACM SIGCOMM*, 2013.
- [25] Open Networking Foundation, "OpenFlow-enabled mobile and wireless networks," *ONF Solution Brief*, 2013.
- [26] Metro Ethernet Forum, "Mobile backhaul - phase 2 implementation agreement," *MEF 22.1*, 2012.
- [27] Metro Ethernet Forum, "Carrier ethernet class of service - phase 2 implementation agreement," *MEF 23.1*, 2012.
- [28] X. Jin, H. H. Liu, R. Gandhi, S. Kandula, R. Mahajan, J. Rexford, R. Wattenhofer, and M. Zhang, "Dionysus: Dynamic scheduling of network updates," in *ACM SIGCOMM*, 2014.
- [29] T. G. Edwards and W. Belkin, "Using SDN to facilitate precisely timed actions on real-time data streams," in *ACM SIGCOMM Workshop on Hot topics in Software Defined Networks (HotSDN)*, 2014.
- [30] T. Mizrahi, O. Rottenstreich, and Y. Moses, "TimeFlip: Scheduling network updates with timestamp-based TCAM ranges," in *IEEE INFOCOM*, 2015.
- [31] T. Mizrahi, E. Saat, and Y. Moses, "Timed consistent network updates," in *ACM SIGCOMM Symposium on SDN Research (SOSR)*, 2015.
- [32] G. R. Ash, "Use of a trunk status map for real-time DNHR," in *International TeleTraffic Congress (ITC-11)*, 1985.
- [33] S. Kandula, I. Menache, R. Schwartz, and S. R. Babbula, "Calendar for wide area networks," in *ACM SIGCOMM*, 2014.
- [34] P. François and O. Bonaventure, "Avoiding transient loops during the convergence of link-state routing protocols," *IEEE/ACM Transactions on Networking (TON)*, vol. 15, no. 6, pp. 1280–1292, 2007.
- [35] L. Vanbever, S. Vissicchio, C. Pelsner, P. François, and O. Bonaventure, "Seamless network-wide IGP migrations," in *ACM SIGCOMM Computer Communication Review*, vol. 41, pp. 314–325, 2011.
- [36] H. H. Liu, X. Wu, M. Zhang, L. Yuan, R. Wattenhofer, and D. Maltz, "zUpdate: updating data center networks with zero loss," in *ACM SIGCOMM*, pp. 411–422, 2013.
- [37] M. Reitblatt, N. Foster, J. Rexford, C. Schlesinger, and D. Walker, "Abstractions for network update," in *ACM SIGCOMM*, 2012.
- [38] N. Pippenger, "On rearrangeable and non-blocking switching networks," *Journal of Computer and System Sciences*, vol. 17, no. 2, pp. 145–162, 1978.
- [39] T. Mizrahi and Y. Moses, "Time-based updates in software defined networks," in *ACM SIGCOMM Workshop on Hot topics in Software Defined Networks (HotSDN)*, 2013.
- [40] T. Mizrahi and Y. Moses, "On the necessity of time-based updates in SDN," in *Open Networking Summit (ONS)*, 2014.
- [41] Open Networking Foundation, "Openflow switch specification," *Version 1.5.0*, 2015.
- [42] Open Networking Foundation, "Openflow extensions 1.3.x package 2," 2015.
- [43] "TIME4 source code," <https://github.com/TimedSDN>, 2014.
- [44] J. M. Kleinberg, "Single-source unsplittable flow," in *37th Annual Symposium on Foundations of Computer Science, FOCS '96, Burlington, Vermont, USA, 14-16 October, 1996*, pp. 68–77, 1996.
- [45] T. Mizrahi and Y. Moses, "Time Capability in NETCONF," draft-mm-netconf-time-capability, work in progress, IETF, 2014.
- [46] T. Mizrahi and Y. Moses, "Using REVERSEPTP to distribute time in software defined networks," in *International IEEE Symposium on Precision Clock Synchronization for Measurement Control and Communication (ISPCS)*, 2014.
- [47] T. Mizrahi and Y. Moses, "Time-based Updates in OpenFlow: A Proposed Extension to the OpenFlow Protocol," technical report, CCIT Report #835, July 2013, EE Pub No. 1792, Technion – Israel Institute of Technology, <http://tx.technion.ac.il/%7Eedw/OFTimeTR.pdf>, 2013.
- [48] T. Mizrahi and Y. Moses, "REVERSEPTP: A software defined networking approach to clock synchronization," in *ACM SIGCOMM Workshop on Hot topics in Software Defined Networks (HotSDN)*, 2014.
- [49] "CPqD OFSoftswitch," <https://github.com/CPqD/ofsoftswitch13>, 2014.
- [50] "Precision Time Protocol daemon," version 2.3.0, <http://ptpd.sourceforge.net/>, 2013.
- [51] "The DeterLab project," <http://deter-project.org/>, 2015.
- [52] B. Lantz, B. Heller, and N. McKeown, "A network in a laptop: rapid prototyping for software-defined networks," in *HotNets*, 2010.
- [53] A. Tavakoli, M. Casado, T. Koponen, and S. Shenker, "Applying NOX to the datacenter," in *HotNets*, 2009.
- [54] A. Tootoonchian, S. Gorbunov, Y. Ganjali, M. Casado, and R. Sherwood, "On controller performance in software-defined networks," in *USENIX Workshop on Hot Topics in Management of Internet, Cloud, and Enterprise Networks and Services (Hot-ICE)*, 2012.
- [55] C. Rotos, N. Sarrar, S. Uhlig, R. Sherwood, and A. W. Moore, "OFLOPS: An open framework for openflow switch evaluation," in *Passive and Active Measurement*, 2012.
- [56] SMPTE Standard 2022-6, "Transport of High Bit Rate Media Signals Over IP Networks (HBRMT)," 2012.
- [57] "Iperf - The TCP/UDP Bandwidth Measurement Tool," <https://iperf.fr/>, 2014.
- [58] ITU-T Y.1563, "Ethernet frame transfer and availability performance," *ITU-T*, 2009.
- [59] D. Arnold and H. Gerstung, "Enterprise Profile for the Precision Time Protocol With Mixed Multicast and Unicast Messages," draft-ietf-tictoc-ptp-enterprise-profile-05, work in progress, IETF, 2015.

APPENDIX A

A TIME EXTENSION TO THE OPENFLOW PROTOCOL

A. Introduction

This section defines a time extension to the OpenFlow protocol. This extension allows the controller to send OpenFlow commands that include an execution time, indicating to the switch *when* the respective command should be performed.

As specified in [19], a bundle is a sequence of (one or more) OpenFlow modification requests from the controller that is applied as a single OpenFlow operation. The controller uses a commit message to apply the set of requests in the bundle. Consequently, the switch applies all messages in the bundle as a single operation or returns an error.

This extension defines *scheduled bundles*; a bundle commit request may include an *execution time*, specifying *when* the bundle should be committed. A switch that receives a scheduled bundle, commits the bundle as close as possible to the execution time that was specified in the commit message.

This document also defines the bundle features message, allowing the controller to retrieve information about the switch's bundle support, and specifically about its scheduled bundle support.

B. How It Works

1) Overview

This extension allows a bundle operation to be invoked at a scheduled time that is determined by the controller. The time-based bundle procedure is illustrated in Figure 14:

- 1) The controller starts the bundle procedure by sending an `OFPBCT_OPEN_REQUEST`, and receives a reply from the switch.
- 2) The controller then sends a set of N `OFPT_BUNDLE_ADD_MESSAGE` messages, for some $N \geq 1$.
- 3) The controller MAY then send an `OFPBCT_CLOSE_REQUEST`. The close request is optional, and thus the controller may skip this step.
- 4) The controller sends an `OFPBCT_COMMIT_REQUEST`. The `OFPBCT_COMMIT_REQUEST` includes two time-related fields: the time flag and optionally the time property. When the time flag is set, it indicates that this is a *scheduled commit*. A scheduled commit request includes the time property field, which contains the scheduled time at which the switch is expected to apply the bundle.
- 5) After receiving the commit message, the switch applies the bundle at the scheduled time, T_s , and sends a `OFPBCT_COMMIT_REPLY` to the controller.

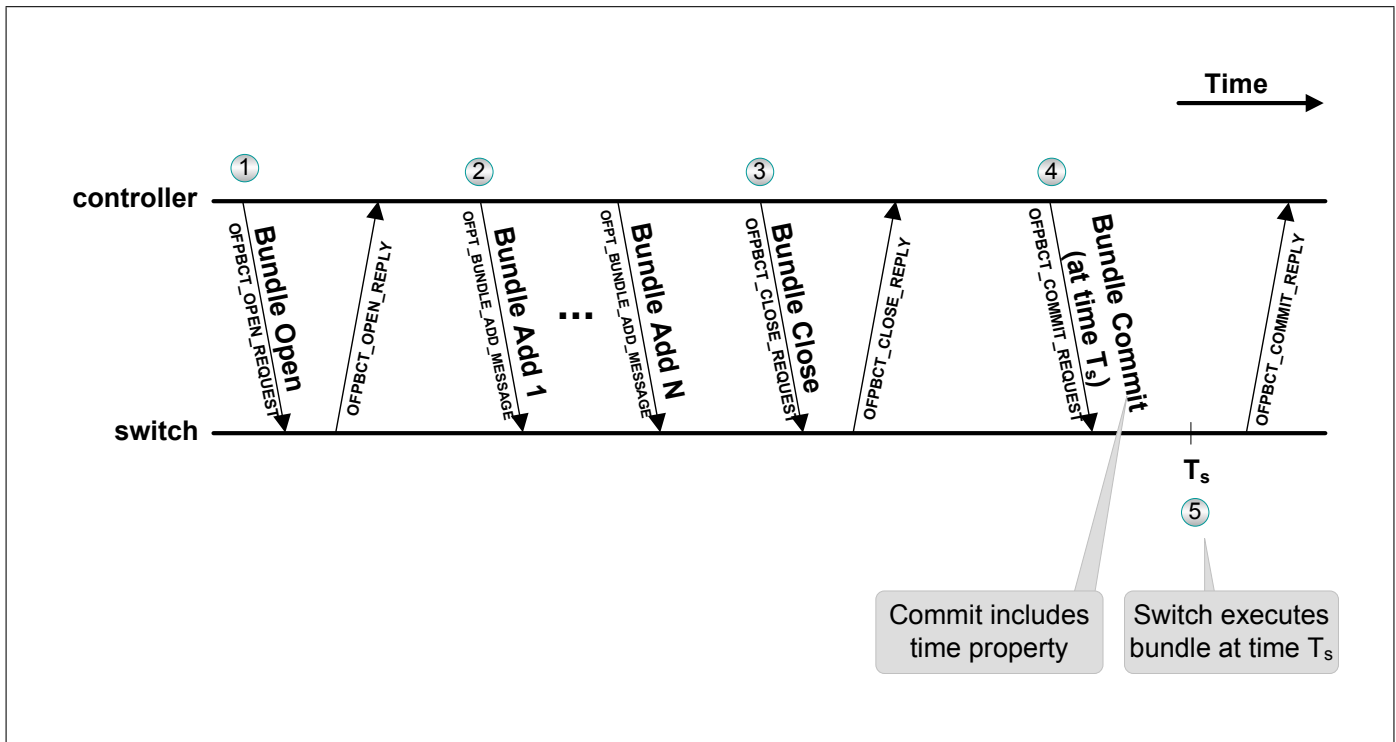


Fig. 14: Scheduled Bundle Procedure

Discarding scheduled bundles. The controller may cancel a scheduled commit by sending an `OFPT_BUNDLE_CONTROL` message with type `OFPBCT_DISCARD_REQUEST`. An example is shown in Figure 15; if the switch is not able to schedule the operation after receiving the commit message, it responds to the controller with an error message (see A-E). This indication may be used for implementing a coordinated update where either all the switches successfully schedule the operation, or the bundle is discarded; when a controller receives a scheduling error message from one of the switches it can send a discard message (step 5' in in Figure 15) to other switches that need to commit a bundle at the same time, and abort the bundle.

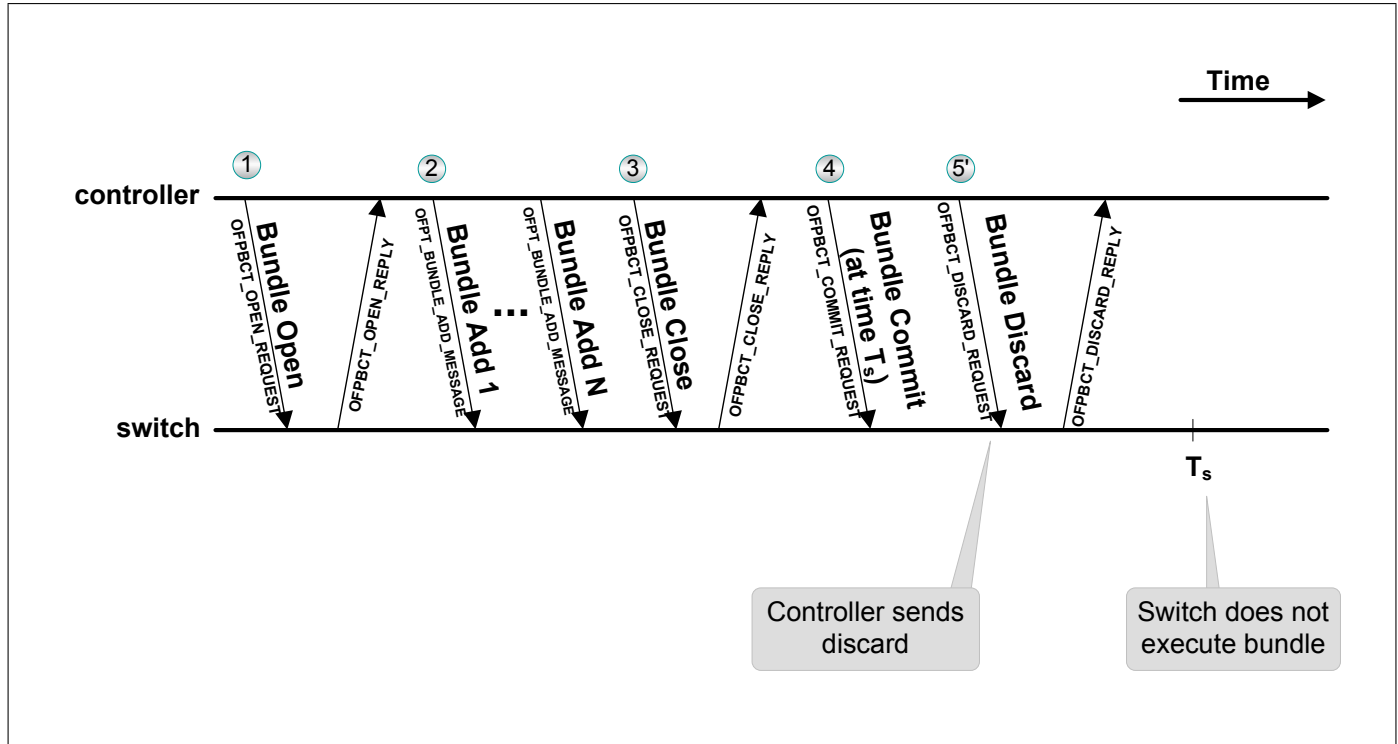


Fig. 15: Discarding a Scheduled Commit

2) Timekeeping and Synchronization

Every switch that supports scheduled bundles must maintain a clock. It is assumed that clocks are synchronized by a method that is outside the scope of this document, e.g., the Network Time Protocol (NTP) or the Precision Time Protocol (PTP).

Two factors affect how accurately a switch can commit a scheduled bundle; one factor is the accuracy of the clock synchronization method used to synchronize the switches' clocks, and the second factor is the switch's ability to execute real-time operations, which greatly depends on how it is implemented.

This document does not define any requirements pertaining to the degree of accuracy of performing scheduled operations. However, every switch that supports the time extension is able to report its estimated scheduling accuracy to the controller. The controller can retrieve this information from the switch using the bundle features message, defined in Section A-D.

Since a switch does not perform configuration changes instantaneously, the processing time of required operations should not be overlooked; in the context of the extension described in this paper the scheduled time and execution time always refer to the start time of the relevant operation.

3) Scheduling Tolerance

When a switch receives a scheduled commit message, it **MUST** verify that the scheduled time, T_s , is not too far in the past or in the future. As illustrated in Figure 16, the switch verifies that T_s is within the *scheduling tolerance* range.

The lower bound on T_s verifies the freshness of the packet so as to avoid acting upon old and possibly irrelevant messages. Similarly, the upper bound on T_s guarantees that the switch does not take a long-term commitment to execute an action that may become obsolete by the time it is scheduled to be invoked.

The scheduling tolerance is determined by two parameters, `sched_max_future` and `sched_max_past`. The default value of these two parameters is 1 second. The controller **MAY** set these fields to a different value using the bundle features request, as described in Section A-D.

If the scheduled time, T_s , is within the scheduling tolerance range, the scheduled commit is performed; if T_s occurs in the past and within the scheduling tolerance, the switch applies the bundle as soon as possible. If T_s is a future time, the switch

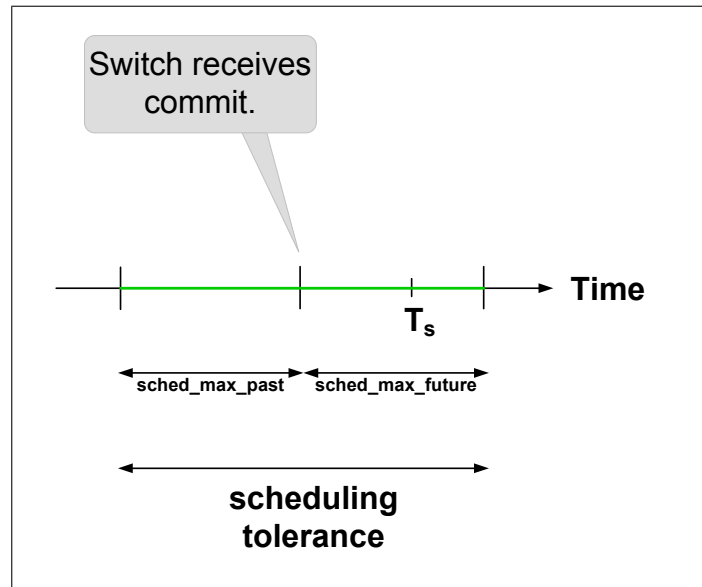


Fig. 16: Scheduling Tolerance

applies the bundle at T_s . If T_s is not within the scheduling tolerance range, the switch responds to the controller with an error message.

C. Time-based Bundle Messages

This section updates Section 7.3.9 of [19]. The reader is assumed to be familiar with Sections 6.8 and 7.3.9 of [19].

The time extension allows bundle commit messages to include a time property, defining when the bundle should be executed.

The time extension defines two time-related fields in OFPBCT_COMMIT_REQUEST messages:

- The time flag, denoted OFPBF_TIME.
- The time property.

All OFPT_BUNDLE_CONTROL messages include the OFPBF_TIME flag. In control messages with type OFPBCT_COMMIT_REQUEST the time flag MAY be set, indicating that the time property field is present. The time property incorporates the time at which the switch is scheduled to apply the bundle.

Control messages with a type that is not OFPBCT_COMMIT_REQUEST MUST have the OFPBF_TIME flag disabled, and this flag is ignored by the switch in these messages.

1) The Time Flag

This document updates ofp_bundle_flags by adding the OFPBF_TIME flag, as follows:

```
/* Bundle configuration flags. */
enum ofp_bundle_flags {
    OFPBF_TIME = 1 << 2, /* Execute in a specific time. */
};
```

2) The Bundle Time Property

This document defines a new bundle property, the time property.

```
/* Bundle property */
struct ofp_bundle_prop_time {
    uint16_t type; /* OFPBPT_TIME */
    uint16_t length; /* Length in bytes = 24 */
    uint8_t pad[4];

    struct ofp_time scheduled_time; /* The scheduled time at which the switch should apply the bundle. */
};
OFP_ASSERT(sizeof(struct ofp_bundle_prop_time) == 24);
```

The type field in the time property is set to the value OFPBPT_TIME, defined as follows:

```
/* Bundle property types. */
enum ofp_bundle_prop_type {
    OFPBPT_TIME = 1, /* Time property. */
};
```

```
};
```

3) Time Format

The time format defined in this extension is based on the one defined in [5]. It consists of two sub-fields; a `seconds` field, representing the integer portion of time in seconds¹⁰, and a `nanoseconds` field, representing the fractional portion of time in nanoseconds, i.e., $0 \leq \text{nanoseconds} \leq (10^9 - 1)$.

```
/* Time format */
struct ofp_time {
    uint64_t seconds;
    uint32_t nanoseconds;
    uint8_t pad[4];
};
OFP_ASSERT(sizeof(struct ofp_time) == 16);
```

As defined in [5], time is measured according to the International Atomic Time (TAI) timescale. The epoch is defined as 1 January 1970 00:00:00 TAI.

D. Bundle Features Request

The bundle features request defined in this document allows a controller to query a switch about its bundle capabilities, including its scheduled bundle capabilities.

This section extends Section 7.3.5 of [19]. The reader is assumed to be familiar with Section 7.3.5 of [19].

The bundle features request is a new multipart message type, the `OFPMP_BUNDLE_FEATURES` message. This document updates `ofp_multipart_type` by adding the `OFPMP_BUNDLE_FEATURES` type, as follows:

```
enum ofp_multipart_type {
    /* Bundle features.
     * The request body is ofp_bundle_features_request.
     * The reply body is struct ofp_bundle_features. */
    OFPMP_BUNDLE_FEATURES = 17,
};
```

1) Bundle Features Request Message Format

The body of the bundle features request message is defined by struct `ofp_bundle_features_request`, as follows:

```
/* Body of OFPMP_BUNDLE_FEATURES request. */
struct ofp_bundle_features_request {
    uint32_t feature_request_flags; /* Bitmap of "ofp_bundle_feature_flags". */
    uint8_t pad[4];

    /* Bundle features property list - 0 or more. */
    struct ofp_bundle_features_prop_header properties[0];
};
OFP_ASSERT(sizeof(struct ofp_bundle_features) == 8);
```

The body consists of a `flags` field, followed by zero or more property TLV fields. The `flags` field, `feature_request_flags`, is defined as follows:

```
/* Flags used in a OFPMP_BUNDLE_FEATURES request. */
enum ofp_bundle_feature_flags {
    OFPBF_TIMESTAMP = 1 << 0, /* When enabled, the current request includes a timestamp, using
                               * the time property */
    OFPBF_TIME_SET_SCHED = 1 << 1, /* When enabled, the current request includes the sched_max_future
                                   * and sched_max_past parameters, using the time property */
};
```

If at least one of the flags `OFPBF_TIMESTAMP` or `OFPBF_TIME_SET_SCHED` is set, the bundle features request includes a time property.

The bundle features properties are specified below.

2) Bundle Features Reply Message Format

If the features request is processed successfully by the switch, it sends a reply to the controller. The body of the bundle features reply message is struct `ofp_bundle_features`, as follows:

```
/* Body of reply to OFPMP_BUNDLE_FEATURES request. */
struct ofp_bundle_features {
    uint16_t capabilities; /* Bitmap of "ofp_bundle_flags". */
};
```

¹⁰The seconds field in IEEE 1588 is 48 bits long. The seconds field used in this extension is a 64-bit field, but it has the same semantics as the seconds field in the IEEE 1588 time format.


```
uint8_t pad[6];

/* Bundle features property list - 0 or more. */
struct ofp_bundle_features_prop_header properties[0];
};
OFP_ASSERT(sizeof(struct ofp_bundle_features) == 8);
```

3) Bundle Features Properties

The optional property fields are defined as TLVs with a common header format, as follows:

```
/* Common header for all bundle feature Properties */
struct ofp_bundle_features_prop_header {
    uint16_t type; /* One of OFPTMPBF_*. */
    uint16_t length; /* Length in bytes of this property. */
};
OFP_ASSERT(sizeof(struct ofp_bundle_features_prop_header) == 4);
```

The currently defined types are as follows:

```
/* Bundle features property types. */
enum ofp_bundle_features_prop_type {
    OFPTMPBF_TIME_CAPABILITY = 0x1, /* Time feature property. */
    OFPTMPBF_EXPERIMENTER = 0xFFFF, /* Experimenter property. */
};
```

The Bundle Features Time Property.

A bundle feature request in which at least one of the flags `OFPBF_TIMESTAMP` or `OFPBF_TIME_SET_SCHED` is set, incorporates the time property. A bundle feature reply that has the `OFPBF_TIME` flag set incorporates the time property.

The time property is defined as follows:

```
struct ofp_bundle_features_prop_time {
    uint16_t type; /* OFPTMPBF_TIME_CAPABILITY. */
    uint16_t length; /* Length in bytes of this property. */
    uint8_t pad[4];

    struct ofp_time sched_accuracy; /* The scheduling accuracy, i.e., how accurately the switch can
    * perform a scheduled commit. This field is used only in bundle
    * features replies, and is ignored in bundle features requests. */
    struct ofp_time sched_max_future; /* The maximal difference between the
    * scheduling time and the current time. */
    struct ofp_time sched_max_past; /* If the scheduling time occurs in the past, defines the maximal
    * difference between the current time and the scheduling time. */
    struct ofp_time timestamp; /* Indicates the time during the transmission of this message. */
};
OFP_ASSERT(sizeof(struct ofp_bundle_features_prop_time) == 72);
```

The time property in a bundle features request includes:

- `sched_accuracy`: this field is relevant only to bundle features replies, and the switch must ignore this field in a bundle features request.
- `sched_max_future` and `sched_max_past`: a switch that receives a bundle features request with `OFPBF_TIME_SET_SCHED` set, should attempt to change its scheduling tolerance values according to the `sched_max_future` and `sched_max_past` values from the time property. If the switch does not successfully update its scheduling tolerance values, it replies with an error message.
- `timestamp`, indicating the controller's time during the transmission of this message. A switch that receives a bundle features request with `OFPBF_TIMESTAMP` set, may use the received timestamp to roughly estimate the offset between its clock and the controller's clock.

The time property in a bundle features reply includes:

- `sched_accuracy`, indicating the estimated scheduling accuracy of the switch. For example, if the value of `sched_accuracy` is 1000000 nanoseconds (1 ms), it means that when the switch receives a bundle commit scheduled to time T_s , the commit will in practice be invoked at $T_s \pm 1\text{ ms}$. The factors that affect the scheduling accuracy are discussed in Section A-B.
- `sched_max_future` and `sched_max_past`, containing the scheduling tolerance values of the switch. If the corresponding bundle features request has the `OFPBF_SET_TIME_TOLERANCE` flag enabled, these two fields are identical to the ones sent by the controller in the request.
- `timestamp`, indicating the switch's time during the transmission of this feature reply. Every bundle feature reply that includes the time property also includes a timestamp. The timestamp may be used by the controller to get a rough estimate of whether the switch's clock is synchronized to the controller's.

E. Errors

As defined in Section 7.5.4 of [19] the switch can send an error message to the controller, which includes a type and a code. This document extends Section 7.5.4 with additional codes, as specified below.

1) Bundle Error

When the switch has an error related to the bundle operation, it sends an error message with type `OFPET_BUNDLE_FAILED`. This document defines the following new codes:

- `OFPBFC_SCHED_NOT_SUPPORTED` - this code is used when the switch does not support scheduled bundle execution and receives a commit message with the `OFPBF_TIME` flag set.
- `OFPBFC_SCHED_FUTURE` - used when the switch receives a scheduled commit message and the scheduling time exceeds the `sched_max_future` (see Section A-B).
- `OFPBFC_SCHED_PAST` - used when the switch receives a scheduled commit message and the scheduling time exceeds the `sched_max_past` (see Section A-B).

The `ofp_bundle_failed_code` is updated as follows:

```
enum ofp_bundle_failed_code {
    OFPBFC_SCHED_NOT_SUPPORTED = 16, /* Scheduled commit was received and scheduling is not supported. */
    OFPBFC_SCHED_FUTURE = 17, /* Scheduled commit time exceeds upper bound. */
    OFPBFC_SCHED_PAST = 18, /* Scheduled commit time exceeds lower bound. */
};
```

2) Bundle Features Error

When the switch has an error related to the `OFPMP_BUNDLE_FEATURES` request, it replies with an error message of type `OFPET_BAD_REQUEST`. The code `OFPBRC_MULTIPART_BAD_SCHED` indicates that the request had the `OFPBF_SET_TIME_TOLERANCE` flag enabled, and the switch failed to update the scheduling tolerance values.

The `ofp_bad_request_code` is updated as follows:

```
enum ofp_bad_request_code {
    OFPBRC_MULTIPART_BAD_SCHED = 16, /* Switch received a OFPMP_BUNDLE_FEATURES request and failed
                                         * to update the scheduling tolerance. */
};
```