# The Case for Data Plane Timestamping in SDN

## Technical Report◇, February 2016

Tal Mizrahi, Yoram Moses*

Technion — Israel Institute of Technology

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

*Abstract*—This paper presents the case for Data Plane Timestamping (DPT). We argue that in the unique environment of Software-Defined Networks (SDN), attaching a timestamp to the header of all packets is a powerful feature that can be leveraged by various diverse SDN applications. We analyze three key use cases that demonstrate the advantages of using DPT, and show that SDN applications can benefit even from using as little as one bit for the timestamp field.

## I. INTRODUCTION

### A. DPT in a Nutshell

Time and synchronized clocks are used in network protocols for various purposes, such as network telemetry [2], [3], [4], Time-Sensitive Networking (TSN) [5], and time-triggered network updates [6]. What if we had a **timestamp** attached to **every** packet in the network?

In this paper we make the case for Data Plane Timestamping (DPT). We argue that adding a timestamp to the header of **every packet** is a powerful tool that is useful for various diverse applications.

DPT is especially relevant in Software-Defined Networks (SDN) for two main reasons: (i) An SDN is a locally-administered environment, where a DPT header can be inserted by **ingress** switches and removed by **egress** switches, and (ii) The current trend in SDN [7], [8], [9] assumes protocol-independent forwarding, providing the flexibility to add and remove any header, and to take (match) decisions based on any header field.

The **Data Plane Timestamp Header (DPTH)** indicates the time at which the packet enters the network. This labeling method is both flexible and uniform; the generic label can be flexibly used by multiple different SDN applications at the same time, and it allows switches to treat each packet consistently, based on its timestamp value.

This paper focuses on three main use cases that demonstrate the merit of DPT:

- **Network telemetry.** DPT allows to measure and monitor the network delay and packet loss using the timestamp, without the need for any additional metadata to be exchanged between switches.

- **Consistent network updates.** The timestamp field can play the role of a configuration *version tag* [10], thereby reducing the management overhead of consistent updates.
- **Load balancing.** The DPTH allows the use of time-division for forwarding elephant flows over multiple paths, allowing higher throughput than other stateless load balancing methods.

DPT allows switches to take packet processing decisions based on the timestamp field. Since match procedures in SDN switches allow header fields to be partially masked [6], [8], we leverage the work of [11] to define DPT-based **time ranges**. Thus, policies or paths can be restricted to specific timestamp ranges.

DPT raises an inevitable question: is it practical to add a new header to **all packets** in the network? Many of the networks in which SDN is deployed or considered use network overlay protocols, e.g., VXLAN [12], Geneve [13], and NSH [14]. These overlay protocols provide inherent extensibility for adding metadata fields, and are therefore DPT-friendly. Furthermore, as we discuss in this paper, a compact timestamp, sometimes even a single-bit, may be sufficient in some applications; since an SDN is a locally-controlled environment, a single-bit DPTH can be accommodated by an unused field in the packet header.

### B. Contributions

The main contributions of this paper are:

- We make the case for attaching a timestamp to **all packets** in SDNs.
- Three key use cases that can benefit significantly from using DPT are presented and analyzed.
- We analyze a simple case where a one-bit timestamp is attached to all packets, demonstrating the benefits that can be achieved using this small overhead.
- Using experimental evaluation we show the merits of DPT in each of the three use cases.

### C. Related Work

Packet timestamping has been proposed for various purposes, such as measurement [15], [2], [3], [16], [17], and clock synchronization [18], [19]. The current paper presents DPT, an approach that adds a timestamp to **all packets**, thereby allowing various SDN applications to use the timestamp for different purposes.

TCP supports an extension [15] that adds a timestamp to all packets, allowing to compute the round-trip-time [15]. The current paper proposes to attach a timestamp to **all** packets, including non-TCP traffic. Furthermore, our solution does not require the end points to be aware of the DPT, since it is used only within the premises of the SDN network.

The work of [20] suggested to use accurate time to schedule consistent updates, thereby reducing the update duration. In this paper we show that the use of **in-band timestamps** eliminates the need for version tags, and reduces the load on the SDN controller. TimeFlip [11] introduced the use of time ranges in switches' TCAMs using an internally-generated timestamp. The current paper generalizes the scope of this work; we show that TimeFlips can be applied at a network-wide scale using a network-wide timestamp, and that DPT-based TimeFlips can be implemented using the ternary match logic of common open source SDN switches.

## II. An Overview of DPT

### A. Timestamping Everything

We propose to timestamp **all** packets. Every packet that enters the network's administrative domain is timestamped. Every packet that is forwarded through the network undergoes the following steps, corresponding to Fig. 1:

1) The ingress switch attaches a DPTH to all packets. The ingress switch may either be a hardware switch or a virtual (software) switch, i.e., a vSwitch.
2) Packets are forwarded through the network with the DPTH. Switches can use the DPTH in their match-action processing. In an environment that uses Virtual Network Functions (VNF), the DPTH can be used by VNFs as well.
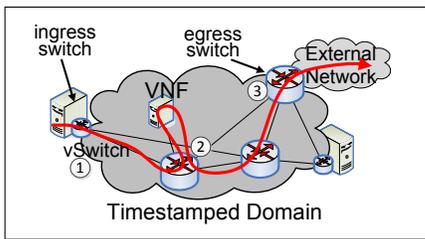3) The egress switch removes the DPTH from the header of every packet.



Fig. 1: Timestamping all packets.

In this paper we focus on the use of DPT in SDN switches. It should be noted that environments that use Network Function Virtualization (NFV) can also benefit from DPT, as the timestamps can be similarly used in the processing of VNFs.

### B. Using the Timestamp

Since a DPTH is integrated into every packet, it can be used in the switches' packet processing procedure. Two of the most promising on-going SDN efforts, P4 [8] and OF-PI [9], allow to flexibly parse and process packet headers.

Specifically the DPTH can be used in the switch match procedure; in addition to conventional match fields such as the 5-tuple, a flow match entry may also include the timestamp field. Since the switch's match process [8], [6] allows to define masks for each match field, it is possible to define time ranges, by masking part of the timestamp field (see Example 1 below). By defining a time range in a flow entry, we are confining the match rule to a specific range of times.

In the context of this paper we focus on two types of time ranges (as defined in [11]): (i) **extremal ranges**, i.e., ranges of the form $T \geq T_0$ (Fig. 2a), and (ii) **periodic ranges**, defining a set of values with a periodic pattern (Fig. 2b).
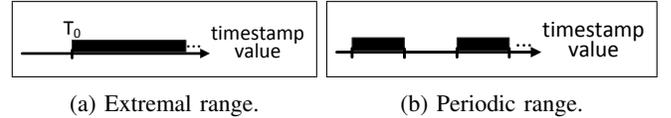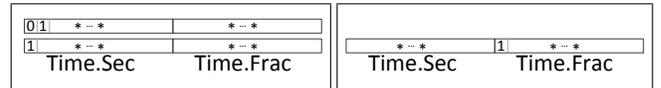


(a) Extremal range.  (b) Periodic range.

Fig. 2: Time ranges.

### C. Timestamp Format

Various different timestamp formats are used in network protocols [18], [19], [2], [21]. In this paper we choose to use the 64-bit NTP timestamp format [18], although the DPT concept applies to other timestamp formats as well.[1] This time format represents the time elapsed since the base date, which is 1 January, 1900. The time format consists of two 32-bit fields: (i) Time.Sec: the integer part of the number of seconds since the base date, and (ii) Time.Frac: the fractional part of the number of seconds.

The size of the DPT timestamp is further discussed in VI.

*Example 1.* Fig. 3 illustrates two time ranges, represented using the NTP time format. The '*' symbol represents bits that are masked.



(a) The extremal range $T \geq 2^{30}$ (b) A periodic range that is active sec. Represented using the $31^{st}$ during the last half of every and $32^{nd}$ bits of Time.Sec. This second. Uses the most significant makes use of two match entries. bit of Time.Frac.

Fig. 3: Time range examples.

As shown in Fig. 3a, representing a time range may require multiple match entries. As discussed in [11], measures can be taken to reduce the number of entries. Specifically, in Sec. IV we show that when a **one-bit** timestamp is used, every time range requires **a single** match entry.

## III. DPT Use Cases

An in-band timestamp can be used for various purposes. We focus on three applications that can greatly benefit from using the timestamp.

---

[1]In fact, DPT can be used with one of the previously defined in-band timestamp encapsulations (e.g., [15], [3], [17]), even though they were defined for different purposes.

## A. Network Telemetry

Performance measurement and monitoring is of key importance in large-scale networks, allowing to detect network faults, anomalies, and congestion, and to enforce a Service Level Agreement (SLA).

*Timestamp-based Measurement:* DPT can be used for accurate delay measurement. As illustrated in Fig. 4, a timestamped packet sent from switch $S_1$ to switch $S_2$ allows $S_2$ to compute the one-way delay from $S_1$ to $S_2$, assuming that the two switches have synchronized clocks. The one-way delay in this case is given by $T_2 - T_1$.
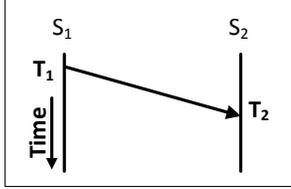


Fig. 4: Timestamp-based delay measurement.

*Color-based Measurement:* The measurement method described above is useful for measuring network **delay**. In contrast, coloring-based passive measurement [22] allows for both **delay** and **loss** measurement.

In the coloring approach of [22] the header of every packet sent between two measurement agents, $S_1$ and $S_2$, includes a binary *color* bit, either '0' or '1'. The color bit divides the traffic into consecutive blocks of packets, allowing $S_1$ and $S_2$ to process each block separately.

According to [22], the color is toggled periodically (e.g., every minute), so that each color is used for a fixed time interval. In the context of SDN this approach would require the SDN controller to periodically update the configuration of $S_1$ each time the color needs to be toggled.

DPT alleviates the need for a color bit. The switches $S_1$ and $S_2$ can derive the color from the in-band timestamp[2] value, by using periodic time ranges, as shown in Fig. 5.
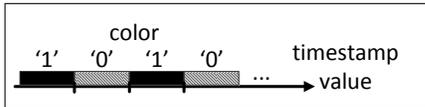


Fig. 5: Coloring based on timestamp ranges.

**The advantage** of our approach is that the SDN controller does not need to periodically perform an update in $S_1$ that toggles the color bit; instead, the color is directly derived from the DPT.

Loss and delay measurement are performed separately for each block.

**Loss measurement (LM).** Each of the two switches that take part in the measurement maintains two counters per flow. The sender $S_1$, maintains $CS0$ and $CS1$, one counter for each color, and the receiver $S_2$, maintains $CR0$ and $CR1$.

---

[2]Note that the DPT is generated by the ingress switch, which may or may not be $S_1$.



(a) Coloring-based loss measurement.
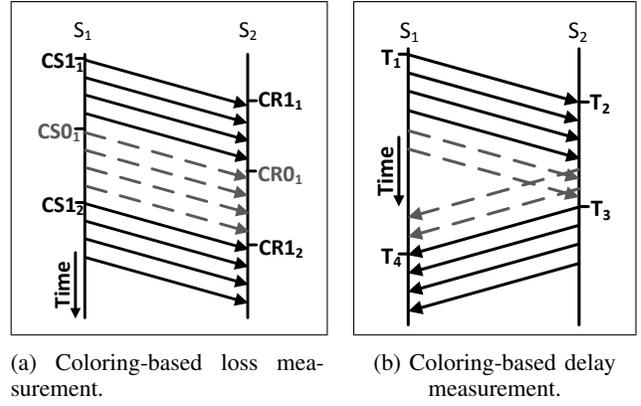
(b) Coloring-based delay measurement.

Fig. 6: Coloring-based method.

Fig. 6a depicts three consecutive blocks of traffic. At the end of each block the controller collects the counter values from each of the switches. Hence, after the second block has completed the controller can compute the $S_1$-to-$S_2$ packet loss during the first block:

$$(CS1_2 - CS1_1) - (CR1_2 - CR1_1) \qquad (1)$$

**Delay measurement (DM).** Each of the two switches uses the first packet of each block as a reference for delay measurement. Switch $S_1$ keeps the time of transmission of the first packet of the block, and $S_2$ keeps the time of reception of this packet. The controller periodically collects these timestamps, to be used for computing the delay.

The one-way delay from $S_1$ to $S_2$, assuming synchronization between the two switches is:

$$T_2 - T_1 \qquad (2)$$

Based on four timestamps collected from two traffic blocks, as shown in Fig. 6b, the controller can compute the two-way delay between the two switches:

$$(T_4 - T_1) - (T_3 - T_2) \qquad (3)$$

Note that the computation of Eq. 3 does not require the two ends to be synchronized.

One of the problems raised in [22] is that packets may arrive to $S_2$ out-of-order, and therefore the timestamp measured by $S_1$ may refer to a different packet than the one measured in $S_2$. DPT provides an inherent solution for this problem; when a switch measures a DM timestamp, it also keeps the DPTH value, extracted from the corresponding packet. The controller can then retrieve the measured timestamps and DPTH values from the two switches, and verify that the two readings have a matching DPTH value.

## B. Consistent Network Updates

Consistent network updates have been widely analyzed in the literature. An update is consistent [10] if every packet is processed by all the switches along its path either according to the 'old' configuration, before the update, or according to the

'new' one. As a test case, we focus on the *two-phase updates* of Reitblatt et al. [10]; all packets include a version tag, indicating for each packet whether it is processed according to the 'old' configuration, or according to the 'new' one. The version tag allows the update to be performed in two phases, as specified in Fig. 7.

---

TWO-PHASE UPDATE

1  Controller sends **new** configuration to switches.
2  Controller enables **new** version tag in ingress switches.

---

Fig. 7: Consistent two-phase update, as in [10].

The DPT approach provides an inherent version tag to all packets. The value of the timestamp progressively increases according to the local clocks of the ingress switches. For each update the controller can define a threshold value $T_{thr}$, such that the **new** configuration is effective during the time range $T \geq T_{thr}$.

---

TIMESTAMP-BASED ONE-PHASE UPDATE

1  Controller sends **new** configuration and $T_{thr}$ to switches. No need to update tagging policy of ingress switches.

---

Fig. 8: Consistent timestamp-based **one-phase** update.

The timestamp allows 'two-phase' updates to be performed using a single[3] phase (Fig. 8). The controller is no longer required to contact the ingress switches to initiate an update. The second phase is replaced by the DPT which automatically assigns timestamps.

Another important advantage of our approach is that it significantly simplifies the management overhead of two-phase updates. For example, consider an SDN application that performs multiple updates. If only a single bit is used for version tagging, then it is not possible to initiate a new update while an update is in progress. By using $2k$ version tag values it is possible to apply $k$ separate updates concurrently under conventional version tagging schemes. The controller would then still need to perform careful bookkeeping of the available version tag values, and after each update the controller would have to wait a sufficient period of time until all the 'old' packets have been flushed from the network before it can reuse the old value. Using DPT all of this overhead can be avoided.

DPT eliminates the need for the SDN application to manage per-flow or per-update version values. Instead, the time-of-day provides global versioning that is guaranteed to be monotonically non-decreasing, while allowing independence between updates.

---

[3]The two-phase approach, as well as our DPT approach, often requires an additional phase for garbage collection, where the old configuration is removed from the switches. This is further discussed in Sec. V.

## C. Timestamp-based Load Balancing

The traditional Equal-Cost Multipath (ECMP) scheme, which balances traffic based on a hash of the packet header, has been in deployment for many years. This method has been shown to be inadequate for elephant flows [23]. Sophisticated methods such as [24] use performance monitoring to dynamically choose the least congested path. Random Packet Spraying (RPS) [25] has been suggested as a simple and stateless alternative for handling elephant flows.

DPT provides a simple mechanism for time-division-based traffic balancing, using periodic time ranges; since the timestamp is embedded in the packet header, we propose to use it as a forwarding criterion. Our approach is simple and stateless, and allows higher performance than existing stateless approaches, as shown in Sec. V.

A simple example is illustrated in Fig. 9a. Traffic from $S_1$ to $S_2$ is forwarded via three paths, A, B, and C. The capacity of each path is depicted in Fig. 9a.
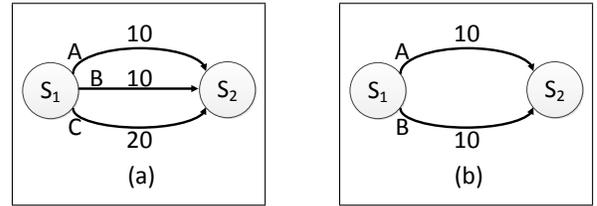


Fig. 9: Load balancing examples.

The DPTH can be used for balancing traffic across the three paths; using periodic time ranges, traffic can be split over the three paths as illustrated in a time-division manner (Fig. 10), providing the desired weight to each path.
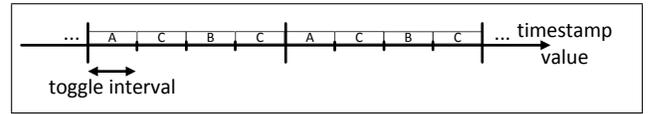


Fig. 10: Load balancing based on periodic timestamp ranges.

## IV. THE POWER OF ONE BIT

In the previous section we presented three applications that can benefit from DPT. One of the main concerns that may arise from using time ranges is its *cost* in terms of the number of timestamp bits added to each packet, and the number of match entries required to represent each time range.

In this section we explore what can be done with a **single-bit** timestamp.

### A. Network Telemetry using One Bit

In this subsection we assume that all data packets are timestamped with a 1-bit timestamp, which is the least significant bit of the Time.Sec field (Sec. II-C). Now the timestamp bit can simply be used as the color indicator, with a toggle interval of one second.

Each of the switches that take part in the measurement uses two separate match entries:

0) Match: timestamp=0
1) Match: timestamp=1

Each entry has its own match counter, yielding a counter for color '0', and a counter for color '1'. The controller reads the counters once per second, allowing to compute the loss rate.

In this example the DPTH provides an inherent color bit without the need for the controller to be involved in the periodic color toggling. Moreover, each switch uses **a single** match rule per color.

### B. Consistent Updates using One Bit

Previous work on dynamic traffic engineering [26], [27] suggested that network paths should be reconfigured periodically, with a period on the order of a few minutes. As an example, we consider a 1-bit timestamp that represents the $7^{th}$ bit of the Time.Sec field, which toggles every 64 seconds.

We propose to use this 1-bit timestamp as a version tag with a periodically toggled value, allowing an SDN traffic engineering application to apply a new set of network paths every 64 seconds.
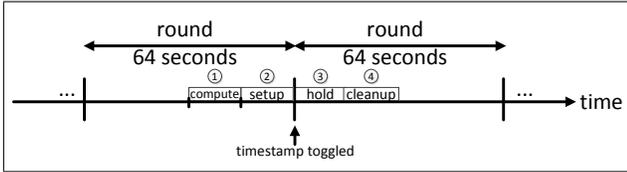


Fig. 11: Periodic consistent updates using a 1-bit timestamp.

Fig. 11 illustrates an example of the periodic schedule of an SDN application that uses this 1-bit timestamp. Each 64-second round is divided into four equal-sized slots: (1) the SDN application re-computes the network paths for the next round, (2) the new configuration is distributed to the switches and installed, (3) after the timestamp bit is toggled, the SDN controller waits until all the old packets have been drained from the network, and (4) the SDN controller removes the configuration of the previous round from the switches.

Each path that is modified between the two rounds requires two match entries during steps (2) and (3), one for the old configuration, and one for the new.

Since the DPTH plays the role of the version tag, the controller does not need to perform the second phase of the two-phase update (Fig. 7), thereby simplifying the 'setup' process compared to conventional two-phase updates.

### C. Load Balancing using One Bit

Consider the network of Fig. 9b. An elephant flow of 20 Gbps needs to be balanced between two 10 Gbps paths. The controller can use timestamp-based rules that define a time division between the two paths.

The total capacity of paths A and B is 20 Gbps. Assuming that packets are typically 1500 bytes long,[4] we roughly have a packet every **0.6 microsecond**. We define a 1-bit timestamp that is the $12^{th}$ bit of the Time.Frac field. The toggle interval of this bit is roughly **0.5 microseconds**. The match rules in switch $S_1$ are defined to be:

0) Match: timestamp=0. Action: output=path A.
1) Match: timestamp=1. Action: output=path B.

Intuitively, we would like the **toggle interval** to be roughly the same as the **inter-packet arrival interval**, as this allows packets to be forwarded in a balanced round-robin-like fashion, without the need for the stateful behavior of round-robin. Hence, configuring the correct toggle interval requires the controller to know the **total bandwidth** of the two paths, and the **MTU**.

## V. EVALUATION

In this section we present experimental evaluation of the three use cases that were presented in previous sections. The goal of these experiments is to quantify the advantages of using DPT in each of the use cases.

The experiments were performed on a testbed of 50 Linux-based machines in the Emulab [28] environment. The experiments included 48 software-based OpenFlow switches, running the open source OFSoftSwitch [29], and Dpctl [30] was used as the controller. We slightly modified the code of OFSoftSwitch, allowing switches to embed the six least significant bytes of the NTP timestamp into the Ethernet source address (eth-src) of each packet. Since the timestamp was piggybacked onto the eth-src field, switches were able to perform match decisions based on this field.
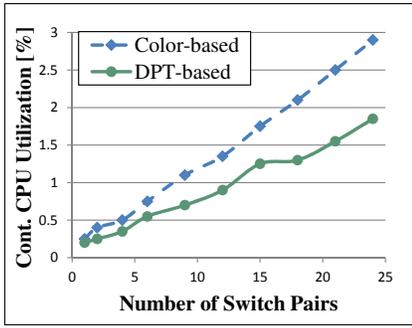
### A. Experiment 1: Telemetry

In this experiment we compared conventional color-based loss measurement [22] to the DPTH-based method of Sec. IV-A.

For each of the two methods we measured the controller's CPU utilization during a measurement of $M$ switch pairs, where each pair of switches performs the measurement as described in Sec. IV-A. We used the Linux 'Top' utility to monitor the average CPU utilization. We used a 1-second measurement interval, and ran each experiment for 100 seconds. We repeated the experiment for various values of $M$.
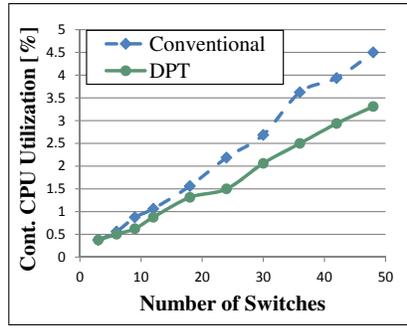
As shown in Fig. 12a, the controller's CPU utilization[5] when using DPT was approximately 30% lower than in the color-based method. The difference is due to the fact that in the color-based method the controller needed to periodically (once per second) read the counters of each of the switches, and **also** to periodically modify the color of each monitored flow. In contrast, in the DPT-based method the controller does not

---

[4]We assume that high-throughput flows are transmitted using maximal-length packets. In this example we assume that the Maximum Transmission Unit (MTU) is 1500 bytes.
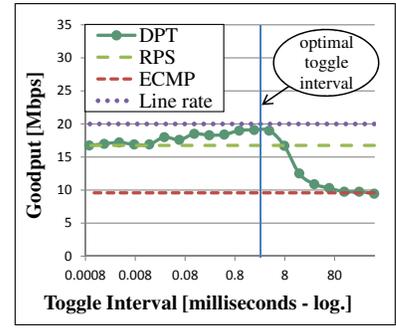
[5]The controller reaches a relatively high CPU utilization, around 3% in an experiment with 48 switches, since the controller we used, Dpctl, is optimized for simplicity at the cost of performance.

(a) Experiment 1: telemetry.



(b) Experiment 2: consistent updates.



(c) Experiment 3: load balancing.
Toggle interval is **chosen** by the SDN
controller based on the link bandwidth,
independent of the flow bandwidth.

Fig. 12: Experimental results.

need to modify the colors of each flow, as each of the switches
derives the color from the DPT.

### B. Experiment 2: Consistent Updates

We compare the controller's load in version-tag-based con-
sistent updates [10] to the DPT-based updates of Sec. IV-B.

We measured the CPU utilization during a periodic consis-
tent update, as described in Sec. IV-B, assuming an update
period of 1 second.[6]

In each round (1 second), we performed a consistent update
that involved $N$ switches. Each update included the following
steps: (i) installing the new configuration in the $N$ switches, (ii)
enabling the new version tag in the ingress switches (this step
is only required in the non-DPT method), and (iii) removing
the old configuration from the $N$ switches. We repeated the
experiment for various values of $N$, up to 48, and each
experiment was performed 100 times. As in [20], we assumed
that $2/3$ of the switches are ingress switches, an assumption
that is applicable to leaf-spine topologies such as Fat-Tree and
Clos. Fig. 12b depicts the CPU utilization in each of the two
methods.

Notably, the DPT-based approach reduces the load on the
controller's CPU by approximately 20% compared to conven-
tional two-phase updates, as the DPT method does not require
step (ii) of the procedure above.

### C. Experiment 3: Load Balancing

We used the topology of Fig. 13, and ran a TCP 'elephant
flow' of 20 Mbps from *src* to *dst* using Iperf [31]. The capacity
of each link (in Mbps) is depicted in the figure.

Switch $S_1$ attached a timestamp to each packet, and $S_2$
performed the balancing between the two 10 Mbps links. In $S_2$
we used one-bit match rules, as described in Sec. IV-C. We
repeated the experiment, each time matching a different bit
in the timestamp field, causing a different **toggle interval**. In
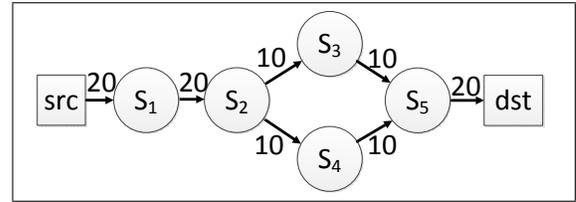each experiment the traffic was run for a period of 10 seconds.



Fig. 13: Load balancing experiment.

Fig. 12c illustrates the measured goodput of the TCP
elephant flow as a function of the toggle interval. The figure
shows the performance of the DPT-based method, compared to
the two stateless approaches discussed in Sec. III-C: Random
Packet Spraying (RPS) [25], and conventional ECMP.

Clearly, ECMP is not optimized for elephant flows, utilizing
a bit under 50% of the network capacity. RPS performs better
than ECMP, but does not reach the full capacity of the links;
since packets are sprayed randomly across the paths, each path
is subjected to an occasional burst of consecutive packets that
were incidentally forwarded to the same path, resulting in
packet drops. These occasional packet drops cause the TCP
algorithm to reduce the traffic rate.

As shown in Fig. 12c, DPT performs significantly better
than the other methods when the controller configures the
toggle interval as discussed in Sec. IV-C, i.e., when the toggle
interval is roughly the same as the inter-packet interval.[7]
Decreasing the toggle interval gracefully reduces the perfor-
mance, asymptotically resembling the performance of RPS.
The reason is that using one of the lower bits of the timestamp
is similar to using a randomized bit. On the other hand, as
the toggle interval is increased beyond a few milliseconds
the performance drops, since a large burst of packets is sent
through a single path in each toggle interval, asymptotically
behaving in an ECMP-like manner.

---

[6]In Sec. IV-B we discussed an update period of 64 seconds. In this
experiment we used a lower update period of 1 second in order to challenge
the controller's performance.

[7]The controller can compute the minimal inter-packet interval based on the
**link** bandwidth, independent of the **flow** bandwidth, which varies over time.

## VI. DISCUSSION

### A. Timestamp Size

We showed that various applications can benefit from DPT, but what is the desired size of the timestamp field? Using a long timestamp field yields costly performance overhead. Sec. IV showed that a one-bit timestamp can suffice for some applications, but that each application may require the timestamp field to represent a different time scale. The size of the timestamp field presents a tradeoff: the most flexible solution, and also the most costly one, is to use a wide timestamp field, such as the format discussed in Sec. II-C. On the other hand, if the timestamp is known to be used for a specific application, a short timestamp field can provide the same functionality and reduce the on-the-wire overhead.

### B. Timestamps vs. Sequence Numbers

The approach we present suggests to include the time-of-day in every packet. Lamport [32] suggested that distributed applications can use *logical clocks*, which produce monotonically increasing *sequence numbers* instead of *timestamps*.

**The advantage of timestamps.** We argue that there are advantages to using a DPT that represents the time-of-day. For instance, in the telemetry use case the DPT defines the traffic blocks to be at fixed time intervals, allowing the controller to sample the counter and timestamp values at fixed time intervals. In the consistent update use case the use of time-of-day allows the controller to plan a schedule of when the update rules need to be installed and removed.

**Clock accuracy.** Using timestamps implies that switches maintain accurately synchronized clocks. Notably, if clocks are not synchronized, then using timestamps is equivalent to using sequence numbers. The required accuracy of the clocks depends on the application that uses DPT. For example, the telemetry application of Sec. IV-A requires an accuracy that is better than 1 second, as the controller retrieves the measurement information from the switches once per second. The consistent update application described in Sec. IV-B will work well even if the clock accuracy is on the order of several seconds.

### C. Impact on Network Protocols

As mentioned in Sec. II-B, the evolving P4 language is already designed in a DPT-friendly way. It is also important to consider how DPT affects the on-the-wire protocols.

**Data plane protocol.** Since the DPTH is designed to be added to every packet, it must be supported by the data plane encapsulation. Fortunately, many of the recently defined data plane encapsulation protocols [12], [13], [14] have flexible support for type-length-value (TLV) or metadata fields in the packet header. Therefore, it is relatively straightforward to define such TLV or metadata fields for the DPT in each of these encapsulations.

**Control plane protocol.** In order to support DPT, two main extensions are required in OpenFlow [6]. Based on the current activity in the Open Networking Foundation (ONF) and in the P4 consortium, we believe that both features are imminent in future versions of OpenFlow.

- The ability to add or remove a timestamp from the packet header. The ability to store a packet's time of transmission or reception is also required, as described in Sec. III-A. We note that the ONF is currently working on OAM support [33], which requires similar timestamping capabilities, and thus these capabilities are likely to be added in future versions of OpenFlow.
- The ability to use the DPTH in match procedures. The ONF's current work on Protocol Independent Forwarding (PIF) [9] will allow the OpenFlow match procedure to be based on any field in the packet header, including the DPTH.

### D. Security Considerations

The security considerations of using DPT are discussed in detail in [17]. In-band timestamping can be used as a means for network reconnaissance. By passively eavesdropping to timestamped traffic, an attacker can gather information about network delays and performance bottlenecks.

The DPT is intended to be used by various diverse applications. Thus, a man-in-the-middle attacker can maliciously modify timestamps in order to attack applications that use the timestamp values.

DPT relies on an underlying time synchronization protocol. Thus, if the time protocol is not properly secured [34], then by attacking the time protocol an attack can potentially compromise the integrity of the DPT.

## VII. CONCLUSION

DPT provides a single header field that can be used for multiple purposes. We have shown three interesting use cases for DPT, and we believe that DPT can be useful for various other SDN applications. Although adding a DPT header to all packets appears to be costly at a first glance, we have shown that the cost of DPT is often reduced to a single bit. Our work suggests that the benefits of DPT in SDNs certainly outweigh the costs.

### REFERENCES

[1] T. Mizrahi and Y. Moses, "The case for data plane timestamping in sdn," in *IEEE INFOCOM Workshop on Software-Driven Flexible and Agile Networking (SWFAN)*, 2016.

[2] ITU-T G.8013/Y.1731, "OAM functions and mechanisms for Ethernet based networks," *ITU-T*, 2015.

[3] C. Kim, P. Bhide, E. Doe, H. Holbrook, A. Ghanwani, D. Daly, M. Hira, and B. Davie, "In-band network telemetry (INT)," technical specification, P4, 2015.

[4] C. Kim, A. Sivaraman, N. Katta, A. Bas, A. Dixit, and L. J. Wobker, "In-band network telemetry via programmable dataplanes," in *ACM SIGCOMM Symposium on SDN Research (SOSR)*, 2015.

[5] IEEE, "Time-Sensitive Networking Task Group," http://www.ieee802.org/1/pages/tsn.html, 2015.

[6] Open Networking Foundation, "OpenFlow switch specification," *Version 1.5.0*, 2015.

[7] P. Bosshart, D. Daly, G. Gibb, M. Izzard, N. McKeown, J. Rexford, C. Schlesinger, D. Talayco, A. Vahdat, G. Varghese, and D. Walker, "P4: Programming protocol-independent packet processors," *ACM SIGCOMM Computer Communication Review*, vol. 44, no. 3, pp. 87–95, 2014.

[8] "The P4 language specification," version 1.0.2, P4 consortium, 2015.

[9] "OF-PI: A Protocol Independent Layer," version 1.1, ONF, 2014.

[10] M. Reitblatt, N. Foster, J. Rexford, C. Schlesinger, and D. Walker, "Abstractions for network update," in *ACM SIGCOMM*, 2012.

[11] T. Mizrahi, O. Rottenstreich, and Y. Moses, "TimeFlip: Scheduling network updates with timestamp-based TCAM ranges," in *IEEE INFOCOM*, 2015.

[12] M. Mahalingam, D. G. Dutt, K. Duda, P. Agarwal, L. Kreeger, T. Sridhar, M. Bursell, and C. Wright, "Virtual extensible local area network (VXLAN): A framework for overlaying virtualized layer 2 networks over layer 3 networks," RFC 7348, IETF, 2014.

[13] J. Gross, T. Sridhar, P. Garg, C. Wright, I. Ganga, P. Agarwal, K. Duda, D. G. Dutt, and J. Hudson, "Geneve: Generic network virtualization encapsulation," draft-ietf-nvo3-geneve, work in progress, IETF, 2015.

[14] P. Quinn and U. Elzur, "Network service header," draft-ietf-sfc-nsh, work in progress, IETF, 2015.

[15] V. Jacobson, B. Braden, and D. Borman, "TCP extensions for high performance," RFC 1323, IETF, 1992.

[16] R. Mittal, N. Dukkipati, E. Blem, H. Wassel, M. Ghobadi, A. Vahdat, Y. Wang, D. Wetherall, D. Zats, *et al.*, "TIMELY: RTT-based congestion control for the datacenter," in *ACM SIGCOMM*, 2015.

[17] R. Browne, A. Chilikin, B. Ryan, T. Mizrahi, and Y. Moses, "Network service header timestamping," draft-browne-sfc-nsh-timestamp, work in progress, IETF, 2015.

[18] D. Mills, J. Martin, J. Burbank, and W. Kasch, "Network time protocol version 4: Protocol and algorithms specification," RFC 5905, IETF, 2010.

[19] IEEE TC 9, "1588 IEEE Standard for a Precision Clock Synchronization Protocol for Networked Measurement and Control Systems Version 2," *IEEE*, 2008.

[20] T. Mizrahi, E. Saat, and Y. Moses, "Timed consistent network updates," in *ACM SIGCOMM Symposium on SDN Research (SOSR)*, 2015.

[21] C. Newman and G. Klyne, "Date and time on the internet: Timestamps," RFC 3339, IETF, 2002.

[22] M. Chen, L. Zheng, G. Mirsky, and G. Fioccola, "IP Flow Performance Measurement Framework," draft-chen-ippm-coloring-based-ipfpm-framework, work in progress, IETF, 2016.

[23] M. Casado and J. Pettit, "Of mice and elephants," *Network Heresy, Nov*, 2013.

[24] M. Alizadeh, T. Edsall, S. Dharmapurikar, R. Vaidyanathan, K. Chu, A. Fingerhut, F. Matus, R. Pan, N. Yadav, G. Varghese, *et al.*, "Conga: Distributed congestion-aware load balancing for datacenters," in *ACM SIGCOMM*, 2014.

[25] A. Dixit, P. Prakash, Y. C. Hu, and R. R. Kompella, "On the impact of packet spraying in data center networks," in *IEEE INFOCOM*, 2013.

[26] 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.

[27] S. Jain *et al.*, "B4: Experience with a globally-deployed software defined WAN," in *ACM SIGCOMM*, 2013.

[28] Emulab — Network Emulation Testbed, http://www.emulab.net, 2015.

[29] "CPqD OFSoftswitch," https://github.com/CPqD/ofsoftswitch13, 2015.

[30] "Dpctl Documentation," https://github.com/CPqD/ofsoftswitch13/wiki/Dpctl-Documentation, 2013.

[31] "Iperf - The TCP/UDP Bandwidth Measurement Tool," https://iperf.fr/, 2014.

[32] L. Lamport, "Time, clocks, and the ordering of events in a distributed system," *Communications of the ACM*, vol. 21, no. 7, pp. 558–565, 1978.

[33] "ONF - Carrier Grade SDN," https://www.opennetworking.org/technical-communities/areas/operator/1909-carrier-grade-sdn, 2015.

[34] T. Mizrahi, "Security requirements of time protocols in packet switched networks," RFC 7384, IETF, 2014.