

# P4Flow: Monitoring Traffic Flows With Programmable Networks

Habib Mostafaei<sup>1b</sup> and Shafi Afridi<sup>1b</sup>

**Abstract**—Cloud providers perform flow monitoring to get insights from the network traffic flows, often using coarse-grained packet counters or packet probing. These approaches give partial information from ongoing flows or introduce significant overhead if the probe packets cross multiple hops with diverse delay and bandwidth to reach the traffic collector. Recently, In-band Network Telemetry (INT) offered by programmable networks, e.g., P4, has gained attention by providing fine-grained network telemetry. Current attempts on INT are inflexible in collecting telemetry for customs flows according to the desired interval. This letter proposes P4FLOW, a flow monitoring tool for cloud provider networks implemented on programmable data planes. P4FLOW allows the providers to monitor a set of desired flows according to their needs. It reduces at least 1.6x the overhead of telemetry packets compared with the existing approaches.

**Index Terms**—Flow monitoring, programmable networks, in-band network telemetry (INT).

## I. INTRODUCTION

THE network infrastructure of the cloud providers often consists of thousands of switches/routers connected through the links [1], [2]. The services, e.g., video streaming or congestion control, running on these infrastructures generate a massive number of traffic flows. The providers need monitoring tools to get insights from ongoing network flows. SNMP, ping, or other legacy tools can provide valuable information on the status of flows. However, this can be improved by employing the recent advances in programmable networks such as P4 [3]. The provider can get insights from the packets/flows and the intermediate devices by using the In-band Network Telemetry (INT) [4] provided by P4. However, getting the flow statistics from each device adds significant overhead to the network. Therefore, the providers desire to collect the flow information from one of the devices in the path from a source to a destination according to the application's need.

Many techniques exist to collect flow statistics [5]–[9]. OpenTM [7] and FlowCover [8] try to minimize the cost of steering the telemetry information by selecting a node from a source to a destination in SDN-based networks. Nevertheless, these approaches add significant overhead to the network.

Manuscript received August 17, 2021; revised August 29, 2021; accepted August 31, 2021. Date of publication September 1, 2021; date of current version November 11, 2021. This work was supported by German Ministry for Education and Research as BIFOLD—Berlin Institute for the Foundations of Learning and Data under Grant 01IS18025A and Grant 01IS18037A. The associate editor coordinating the review of this letter and approving it for publication was V. Eramo. (Corresponding author: Habib Mostafaei.)

The authors are with the Department of Telecommunication Systems, Fachgebiet Internet Architecture and Management, Technische Universität Berlin, 10587 Berlin, Germany (e-mail: habib@inet.tu-berlin.de; safridi@inet.tu-berlin.de).

Digital Object Identifier 10.1109/LCOMM.2021.3109793

The goal of INT-path [6] and P2INT [9] is to cover all the networks' links to collect the flow statistics in a P4-based network. The work in [10] concentrates on coordinating of probe packets to get insights from the network flows. However, these approaches focus on coordinating probe packets to collect the flow statistics without considering the overhead of steering probe packets to the collector.

In this letter, we propose a flow monitoring tool P4FLOW that uses simple additive weighting technique to find the best network devices to report the flow telemetry in a P4 network. It adds less overhead to the network while providing the flexibility to the users to poll the flow telemetry according to their desired intervals and flows. We summarize our contributions as follows.

- We propose a simple additive weighting-based technique to select a leader node to poll the flow statistics.
- The users of P4FLOW can adjust the interval of having flows telemetry information according to the applications' needs.
- We report that P4FLOW injects significantly fewer packets compared with the existing techniques and improves the flow completion time.

We make our code publicly available and our results fully reproducible.<sup>1</sup>

## II. PROBLEM FORMULATION

We model the network as an undirected graph  $G=(P,E)$  in which  $P = \{p_1, p_2, \dots, p_m\}$  is a set of P4 switches and  $E = \{e_1, e_2, \dots, e_m\}$  is the set of edges among them. The network steers the traffic of an active flows  $\mathcal{F} = \{f_1, f_2, \dots, f_n\}$  where each flow  $f$  consists of a set of P4 switches from a source to a destination. We can report the flow information using each P4 switch along the path, but this solution adds additional overhead to the network due to injecting a large number of packets. Therefore, our goal is to select a P4 switch, i.e., a leader switch, along each flow  $f$  to report the flow information to the controller.<sup>2</sup> We call  $\mathcal{L}_f$  as the leader switch for flow  $f$ . We assume that the first node in the graph is the controller of the network being in charge of collecting the flow information. To avoid possible congestion imposed by the flow telemetry packets, we aim to steer the traffic for telemetry packets via less congested links with high available bandwidth.

### A. Delay

In delay-sensitive applications, the flow information should be forwarded to the controller with a short delay to take a

<sup>1</sup><https://github.com/mostafaei/P4Flow>

<sup>2</sup>We use the general term of the controller in this work since in our implementation the P4 runtime plays the role of the controller.

suitable reaction to the network changes. In such scenarios, minimizing the delay between the leader switch and the controller is desirable. Let  $D_{ij}$  denotes the link latency from switch  $p_i$  to  $p_j$ . We denote the sum of link delays along the path from the leader switch  $\mathcal{L}_f$  to the controller for flow  $f$  as  $D(f)$ . Therefore, our goal is to minimize the value of  $D(f)$  for each flow in  $\mathcal{F}$ .

### B. Bandwidth

We desire to send the flow telemetry to the controller using the links with the highest available bandwidth. In that case, the link with the minimum available bandwidth specifies the available bandwidth. Let  $B_{ij}$  denotes the available bandwidth from switch  $p_i$  to  $p_j$ . We can compute the available bandwidth for flow  $f$ , i.e.,  $B(f)$ , by taking the minimum available bandwidth of all links along the path from  $\mathcal{L}_f$  to the controller. Herein, our goal is to maximize  $B(f)$  for each flow in  $\mathcal{F}$ .

### C. Objective Functions

The main goal is to minimize the delay of steering the flow telemetry to the controller while obeying the bandwidth constraints of the links. We use a simple-additive weighting method [11] to transform the multi-objective selection problem into a single one. We define our objective function  $F$  as the weighted sum of normalized values for delay and bandwidth for each flow  $f$  as follows.

$$F(f) = w_d \times \frac{D(f) - D_{min}}{D_{max} - D_{min}} + w_b \times \frac{B(f) - B_{min}}{B_{max} - B_{min}} \quad (1)$$

where  $w_d, w_b \geq 0$  and  $w_d + w_b = 1$ , are weights for link latency and the available bandwidth.  $D_{min}(D_{max})$  and  $B_{min}(B_{max})$  are the minimum (maximum) value for overall expected latency and available bandwidth. We normalize the values of link latency and available bandwidth to adjust the selection metrics. After normalization, each value varies in the range of  $[0, 1]$ , where the value of 0 corresponds to the worst value, i.e.,  $D(f) = D_{max}$  and  $B(f) = B_{min}$ , and the value of 1 corresponds to the best value, i.e.,  $D(f) = D_{min}$  and  $B(f) = B_{max}$ . Therefore, our objective is to minimize  $F(f)$ , mathematically,

$$\min_{\forall f \in \mathcal{F}} F(f)$$

This model can be extended to include other parameters, e.g., the cost of sending the flow information, in selecting  $\mathcal{L}_f$  for each flow  $f$ .

## III. THE ALGORITHM

In this section, we first present a greedy approach to select a switch in a flow to send the flow information. Then, we explain the P4 implementation of P4FLOW in detail.

### A. The Switch Selection Algorithm

The control plane of P4FLOW takes the set of flows  $\mathcal{F}$  and the coefficient for the link parameters, i.e., latency, and bandwidth, as the inputs. Then, it checks the value of the objective function in eq. 1 and selects the switch in the path

with minimum value as the leader switch to send the flow information. Algorithm 1 presents the pseudo code of the leader selection algorithm in P4FLOW.

---

### Algorithm 1: Leader Node Selection

---

```

input : Network graph  $G$ , links delay and bandwidth, flow list ( $\mathcal{F}$ )
output: A list of leader switches  $\mathcal{L}$ 

1 Function findLeaderSwitch():
2    $\mathcal{L} = \emptyset$ ;
3   for each  $f \in \mathcal{F}$  do
4      $t = \infty$ ;
5     for  $p \in f$  do
6        $F_p = \text{apply eq. 1}$ ;
7       if  $F_p < t$  then
8          $t = F_p$ ;
9          $\mathcal{L}_f = p$ ;
10      end
11    end
12     $\mathcal{L} = \mathcal{L} \cup \mathcal{L}_f$ ;
13  end
14  return  $\mathcal{L}$ ;
15 End Function
    
```

---

The algorithm checks the possibility of each switch in the flow for being a leader switch to report the flow information. To do so, the leader selection algorithm checks its distance to the controller using the inter-node delays and their available bandwidth. There might be multiple paths from a switch to the controller, and we use the Dijkstra algorithm to select the best path. It is a single-source shortest path algorithm that works with non-negative edge weights. The running time of Dijkstra's algorithm is less than other algorithms such as Bellman-Ford's one for the same problem. Furthermore, the time complexity of Dijkstra's algorithm can be improved using the Fibonacci heap. The network administrators can adjust each link's cost according to the objective function in eq. 1 since it has a non-negative value. After selection of the best path, we compute the objective function of eq. 1 to check the suitability of the switch. We select the switch with a minimum value of  $F$  in the flow as the leader. Then, the algorithm adds this node to the list of leader switches  $\mathcal{L}$ . We repeat the same procedure for all the flows.

### B. P4 Implementation

P4FLOW can monitor all the packets of all the flows in the network. However, this adds significant overhead to the network, which operators are reluctant to do. Therefore, P4FLOW allows the operator to define a time interval for the flow telemetry. For example, if the flows are sensitive to the delay, the threshold value for the time interval can be chosen as a small value. While for other traffic flows, a large threshold value is desirable to reduce the overhead. This value can be tuned using the control plane rules. We use a P4 register to keep the threshold value for the interval. To do so, the network administrators can set the initial value for the threshold depending on the services offered by the cloud service providers. If the flows carry the traffic of the critical business transaction information, the threshold value for the time interval has to be small to identify the problem quickly [2]. For example, Akami [12] in 2017 report that each 100ms of network delay results in a 7% revenue reduction.

```

register<bit<48>>(1) pkt_thd_reg;
register<bit<48>>(MAX_ENTRIES) timer_reg;
action do_clone_e2e(){
  clone3(CloneType.E2E,E2E_CLONE_SESSION_ID,meta); }
action mark_packet(){
  meta.flag = 1;
  data=hdr.ipv4.srcAddr ++ hdr.ipv4.dstAddr ++ hdr.tcp.
  srcPort ++ hdr.tcp.dstPort ++ hdr.ipv4.protocol ++
  meta.last_timestamp ++ meta.packet_size;
}
table generate_clone{
  actions = {
    do_clone_e2e;
    NoAction; }
  default_action = NoAction();
}
apply {
  pkt_thd_reg.read(meta.pkt_thd, 0);
  time_t timer_cnt;
  time_t timer_tmp;
  time_t cur_time = standard_metadata.
  ingress_global_timestamp;
  timer_reg.read(timer_tmp , (bit<32>) meta.FlowID);
  timer_cnt=cur_time-timer_tmp;
  if(!IS_E2E_CLONE(standard_metadata)){
    if (cur_time>meta.pkt_thd ){
      timer_reg.write((bit<32>)meta.FlowID, timer_cnt);
      mark_packet ();
    }
    if(meta.flag == 1){
      generate_clone.apply(); }
  }
  else{
    //fill telemetry header data
  }
}

```

Fig. 1. Egress control flow of P4FLOW.

For other flows, such as the flows that carry the email data, the time interval for the threshold can be set to a high value.

P4FLOW consists of two main phases, namely, *tracking flows* and *collecting telemetry*. It performs flow tracking in the ingress control flow of each switch to specify where the packet has to be sent, i.e., to the next hop or the controller. Collecting flow telemetry is carried out in the egress control flow. We now explain the P4 implementation of our tool.

1) *Tracking Flows*: To track the flows, we need to keep the state of each flow in the leader switches. Upon entering a packet of a flow to the ingress control flow, P4FLOW performs a matching based on the source and destination IP addresses provided by the control plane and assigns an ID to each flow, i.e., FlowID. Then, it checks if the packet belongs to that flow or it is a telemetry packet. In both cases, P4FLOW forwards the packets according to the forwarding rules either to the destination of the flow or the controller. We assign a unique ID for each flow belonging to the traffic of different flows, while for all the telemetry packets we use one ID to save memory space of the switches.

2) *Collecting Telemetry*: We use P4 registers to keep track of the flow statistics at the egress control flow. Each P4 register can keep the state of the flows in the memory of the switch. We get the timestamp of each packet from `egress_global_timestamp` metadata. Fig. 1 shows the part of the egress control flow implementation of P4FLOW.

In Fig. 1, we define two registers, namely, `pkt_thd_reg` and `timer_reg`, to store the time interval threshold value to report the flow statistics and the timestamp of the packets for

each flow, respectively. We initialize the `timer_reg` to zero to keep track of the flows using their FlowID. Upon receiving a packet of a flow, P4FLOW updates the corresponding register value of `timer_reg` associated with that flow. The algorithm checks this value with a user-defined threshold value to send the flow information. To collect the statistics regarding the number of packets for each flow, we can use P4 Counters or `packet_length` metadata field. P4FLOW sends the flow information to the controller when the timer reaches the threshold value. We now explain the procedure of sending the flow telemetry.

3) *Packet Generation*: There is no built-in mechanism in P4 to generate the packets, but we can use a copy of the current packet to send it to the desired egress port. One example solution is to clone the packet. Therefore, we use the `clone egress` to `egress` or `clone_e2e` primitive to mark the flow information on the packet header similar to [13]. The `mark_packet` action in Fig. 1 concatenates all the telemetry data to one bitstring and appends it to the cloned packet. This data field includes source and destination IP addresses, source and destination ports, and protocol fields to distinguish different flows. Additionally, we concatenate the last timestamp the telemetry updated and the size of transmitted data to this field. We define a one-bit metadata variable called `flag` to specify if the packet should be cloned at the egress control flow by applying `generate_clone` table.

## IV. EVALUATION

In this section, we first report the performance of our system on two custom graphs. Then, we measure the impact of the telemetry packets on the Flow Completion Time (FCT) on a real-network topology taken from [14].

### A. Performance on Custom Graphs

We generate two graphs, namely, Waxman and Erdős-Rényi, each one with 200 nodes to compare the performance of different algorithms. We use  $\alpha = 0.60$  and  $\beta = 0.06$  as the values for the model parameters in generating Waxman graph in Python. Similarly, we set  $p = 0.05$  as the probability of edge creation in Erdős-Rényi graph. Both graphs are widely used by the network community for research purposes. The Waxman graph has 644 edges, while the Erdős-Rényi graph has 1004 edges. We set each link latency value of the graphs by randomly selecting a value from the U.S network latency reported by AT&T [15]. Finally, we randomly assign a bandwidth value in Mbps for each link in the range of 10 to 40 Mbps. We generate 100 to 1000 uniformly random flows in which each flow  $f$  has maximum hops of six and assume that they are assigned using the traffic engineering techniques such as the ones in [2]. Therefore, P4FLOW leverages the spare network capacity to steer the telemetry packets. We measure the reporting delays and the message overhead to the controller/collector. We also report the running times of the algorithms.

We compare the performance of P4FLOW with those of OpenTM [7] and FlowCover [8] as two state-of-the-art flow



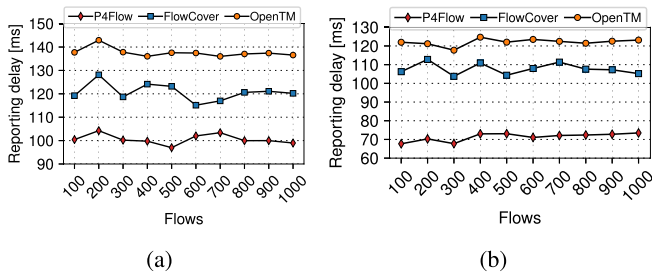


Fig. 2. The reporting delays of the systems on; a) Waxman graph. b) Erdős-Rényi graph.

monitoring tools. The network operators can adopt these approaches to monitor custom flows. However, several other works concentrate on the different aspects of flow monitoring such as coordinating the INT packets or monitoring all links.

OpenTM [7] considers the traffic matrix of the network and aims to minimize the number of telemetry packets. It offers different mechanisms for switch selection such as random, the last switch, and the least-loaded one. We implement the last switch selection in this work. FlowCover [8] aims to minimize the flow monitoring cost with high accuracy. It selects the switch that carries the highest number of flows in the network. To implement FlowCover, we check all switches in all the flows and pick the most frequent ones as the leaders. To have a fair comparison with OpenTM and FlowCover, we use the same P4 code of P4FLOW to run the experiments but with the chosen leader nodes by each approach. Herein, selecting different values for the time interval can negatively impact the obtained results.

1) *Reporting Delay*: Fig. 2 reports the flow statistics message delays for all systems on Waxman and Erdős-Rényi graphs. We take the sum of the link delay of each hop from the leader switches to the controller. The reporting delay of P4FLOW on average is 1.2x and 1.4x less than FlowCover and OpenTM on Waxman graph (see Fig. 2a). This improvement on Erdős-Rényi graph is 1.57x and 1.71x compared with the other approaches since the graph has more edges (see Fig. 2b). The main reason for such an improvement of P4FLOW is due to the parameters, e.g., inter-node delay, that considers in selecting leader switches.

2) *Overhead*: We compute the overhead of the telemetry packets for each flow  $C(f)$  as follows.

$$C(f) = \mathcal{H}_f \times \mathcal{S}, \quad (2)$$

where  $\mathcal{H}_f$  and  $\mathcal{S}$  are the number of hops from the  $\mathcal{L}_f$  to the controller and probe packet size, respectively.

Fig. 3 reports the flow statistics overhead of all systems for the same graphs. We use the same message size in all approaches to have a fair comparison. (see Fig. 3a) shows that the message overhead of P4FLOW on average is 1.63x and 1.78x less than FlowCover and OpenTM on the Waxman graph. We have similar improvement on Erdős-Rényi graph and P4FLOW has 1.60x and 1.76x less overhead compared with the other approaches (see Fig. 3b). The main reason for such an improvement is that P4FLOW uses fewer hops to send the flow statistics messages.

3) *Running Time*: We now measure the running time of the leader selection algorithm in P4FLOW, OpenTM, and

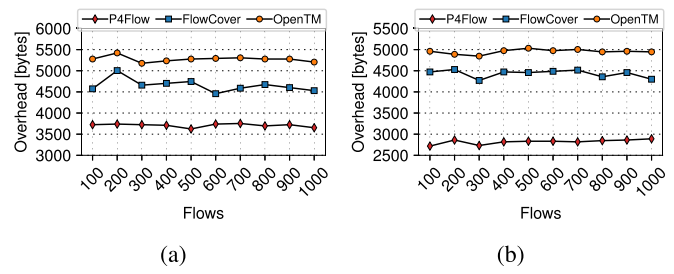


Fig. 3. The message overhead of the systems on; a) Waxman graph. b) Erdős-Rényi graph.

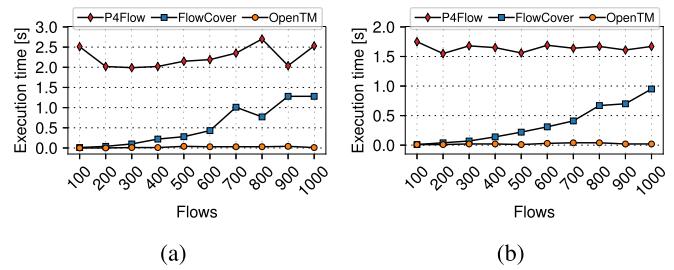


Fig. 4. The running time of the systems on; a) Waxman graph. b) Erdős-Rényi graph.

FlowCover. Fig. 4 shows the running time of all systems for the same graphs. OpenTM selects the last node of a flow to report the flow statistics. Thus, its running time is constant and close to zero. The running time of FlowCover increases by increasing the number of flows because it has to check which node is the most repeated in all flows. P4FLOW first ranks the node and then applies the leader selection algorithm to choose a node based on objective function in eq. 1. Therefore, its running time is higher than the other two approaches. However, this running time has less impact on the performance of the system since the network administrator can proactively execute it.

## B. Performance on WAN Networks

We now report the impact of flow telemetry on the FCT of different systems using an example wide-area network.

1) *Testbed*: We conduct the simulation using Mininet network emulator on a VM with an Intel Xeon CPU AMD Opteron(TM) Processor 6272 2.1GH with 32 GB RAM and 16 CPU cores running Ubuntu server 18.04. Mininet is the defacto emulation tool used by the P4 consortium to test the functionality of the P4 programs. We also apply the recommendation [16] of the P4 consortium to gain high throughput from BMv2 switches in Mininet.

2) *Workloads*: We use one of the most popular empirically derived realistic workloads, i.e., FaceBook [17]. The traffic distribution in this workload is heavy-tailed. We leverage the traffic generator in [18] to generate the desired flows according to a Poisson distribution of the workload and network load. The traffic generator generates traffic from 50 different sources that are randomly chosen from the network topology. The network load varies in the range 10% and 60%. The traffic generator can generate different flow sizes, but we limit them to the small size flows, i.e., flow size < 100 KB due to the

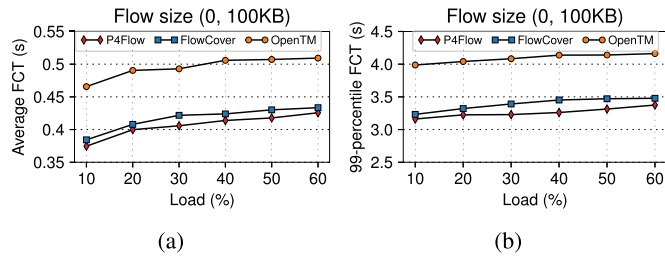


Fig. 5. Comparisons of FCT for P4FLOW, OpenTM, and FlowCover. a) Average FCT of small flows. b) the 99<sup>th</sup> percentile of small flows.

limitations of our links in Mininet. However, this can be tuned according to the capacity of each link. We generate 1000 flows from each source with 0.8 ms as the threshold interval and report the averaged results over 10 different runs.

3) *Topology*: We compare the FCT of P4FLOW with those of OpenTM and FlowCover on RocketFuel [14] topologies such as AS3967 as an example and report some representative results. We use the same links delay and bandwidth for these topologies since their exact values are unavailable.

4) *FCT Comparison*: We now measure the impact of the leader nodes' selection on the FCT of flows. The goal is to understand the impact of the flows telemetry' overhead on FCTs. Fig. 5 shows average FCT and the 99<sup>th</sup> percentile of flows of P4FLOW, FlowCover, and OpenTM on AS3967. P4FLOW improves up to 1.04x and 1.24x the FCT of flows over FlowCover and OpenTM, respectively. Furthermore, this improvement in the 99<sup>th</sup> percentile FCT of flows is up to 1.08x and 1.29x compared with FlowCover and OpenTM. The results confirm that selecting proper leader nodes impacts the FCT of the flows in the network.

## V. PRACTICAL APPLICATION SCENARIOS

In this section, we briefly explain some practical application scenarios that P4FLOW can help. The first application scenario is geo-distributed streaming applications running on distributed processing systems such as Apache Flink. Many organizations leverage these systems to get invaluable insights from their customers by collecting information regarding the users' clicks on the websites or finding trends on social networks [19]. Depending on the queries running on this massive amount of data generated by the users, the network administrators can tune the time interval of collecting flow statistics. The other scenarios are video streaming or online gaming. The network operators can set the threshold value to gather the flow statistics depending on the sensitivity of the traffic flows. For example, delaying the flows in the video streaming scenarios can degrade the quality of experience (QoE) [20].

## VI. CONCLUSION

In this letter, we proposed a network flow monitoring tool for programmable networks. Our tool can reduce the overhead

of reporting the flow statistics by considering the properties of the inter-node links such as delay and bandwidth. The P4-based implementation of our approach showed that it can efficiently report the flow telemetry while improving the flow completion time. Furthermore, it allows the provider to tune collecting the flow statistics according to the applications' needs.

## REFERENCES

- [1] Y. Zhou *et al.*, "Flow event telemetry on programmable data plane," in *Proc. SIGCOMM*, 2020, pp. 76–89.
- [2] H. Mostafaei, M. Shojafar, and M. Conti, "TEL: Low-latency failover traffic engineering in data plane," *IEEE Trans. Netw. Service Manage.*, early access, Jul. 26, 2021, doi: [10.1109/TNSM.2021.3099620](https://doi.org/10.1109/TNSM.2021.3099620).
- [3] P. Bosshart *et al.*, "P4: Programming protocol-independent packet processors," *SIGCOMM Comput. Commun. Rev.*, vol. 44, no. 3, pp. 87–95, Jul. 2014.
- [4] The P4.org Applications Working Group. (2020). *In-Band Network Telemetry (INT) Dataplane Specification V2.1*. [Online]. Available: <https://github.com/p4lang/p4-applications/tree/master/docs>
- [5] J. Kučera, D. A. Popescu, H. Wang, A. Moore, J. Kořenek, and G. Antichi, "Enabling event-triggered data plane monitoring," in *Proc. Symp. SDN Res. (SOSR)*, Mar. 2020, pp. 14–26.
- [6] T. Pan *et al.*, "INT-path: Towards optimal path planning for in-band network-wide telemetry," in *Proc. IEEE Conf. Comput. Commun. (INFOCOM)*, Apr. 2019, pp. 487–495.
- [7] A. Tootoonchian, M. Ghobadi, and Y. Ganjali, "OpenTM: Traffic matrix estimator for OpenFlow networks," in *Passive and Active Measurement, A. Krishnamurthy and B. Plattner*, Eds. Berlin, Germany: Springer, 2010, pp. 201–210.
- [8] Z. Su, T. Wang, Y. Xia, and M. Hamdi, "FlowCover: Low-cost flow monitoring scheme in software defined networks," in *Proc. IEEE Global Commun. Conf.*, Dec. 2014, pp. 1956–1961.
- [9] A. G. Castro *et al.*, "Near-optimal probing planning for in-band network telemetry," *IEEE Commun. Lett.*, vol. 25, no. 5, pp. 1630–1634, May 2021.
- [10] R. Hohemberger *et al.*, "Orchestrating in-band data plane telemetry with machine learning," *IEEE Commun. Lett.*, vol. 23, no. 12, pp. 2247–2251, Dec. 2019.
- [11] E. Triantaphyllou, *Multi-Criteria Decision Making Methods*. New York, NY, USA: Springer, 2000.
- [12] J. Young and T. Barth. (2017). *Web Performance Analytics Show Even 100-Millisecond Delays Can Impact Customer Engagement and Online Revenue*. [Online]. Available: <https://bit.ly/3xzNkyT>
- [13] R. Joshi, T. Qu, M. C. Chan, B. Leong, and B. T. Loo, "BurstRadar: Practical real-time microburst monitoring for datacenter networks," in *Proc. 9th Asia-Pacific Workshop Syst. (APSys)*, Aug. 2018, pp. 1–8.
- [14] N. Spring, R. Mahajan, and D. Wetherall, "Measuring ISP topologies with Rocketfuel," *ACM SIGCOMM Comput. Commun. Rev.*, vol. 32, no. 4, pp. 133–145, 2002.
- [15] ATT Global IP Network Latency. *U.S. Network Latency*. Accessed: Jun. 25, 2021. [Online]. Available: [https://ipnetwork.bgtmo.ip.att.net/pws/network\\_delay.html](https://ipnetwork.bgtmo.ip.att.net/pws/network_delay.html)
- [16] P4 Language Consortium. (2020). *Performance of BMV2*. [Online]. Available: <https://github.com/p4lang/behavioral-model/blob/main/docs/performance.md>
- [17] A. Roy, H. Zeng, J. Bagga, G. Porter, and A. C. Snoeren, "Inside the social network's (datacenter) network," in *Proc. ACM Conf. Special Interest Group Data Commun. (SIGCOMM)*, Aug. 2015, pp. 123–137.
- [18] W. Bai, L. Chen, K. Chen, and H. Wu, "Enabling {ECN} in multi-service multi-queue data centers," in *Proc. 13th USENIX Symp. Netw. Syst. Design Implement. (NSDI)*, 2016, pp. 537–549.
- [19] H. Mostafaei, S. Afridi, and J. H. Abawajy, "SNR: Network-aware geo-distributed stream analytics," in *Proc. IEEE/ACM 21st Int. Symp. Cluster, Cloud Internet Comput. (CCGrid)*, May 2021, pp. 820–827.
- [20] H. Mao, R. Netravali, and M. Alizadeh, "Neural adaptive video streaming with pensieve," in *Proc. Conf. ACM Special Interest Group Data Commun. (SIGCOMM)*, Aug. 2017, pp. 197–210.