

Virtual TCAM for Data Center Switches

Qasim Maqbool
xFlow Research Inc.
qasim.maqbool@xflowresearch.com

Junaid Zulfiqar
xFlow Research Inc.
junaid.zulfiqar@xflowresearch.com

Sohaib Ayub
Lahore University of Management Sciences
13030022@lums.edu.pk

Aamir Shafi
xFlow Research Inc., National University of Sciences and
Technology, Pakistan
aamir.shafi@xflowresearch.com
aamir.shafi@seecs.edu.pk

Abstract— Software Defined Data Center Networks utilize a large number of rules to define routes for network flows. SDN controllers install rules at each switch in order to maximize metrics such as efficiency, delay or energy consumption. This requires each switch to store a substantial number of rules, which is not possible for current Top-of-Rack (ToR) switches as they usually have capacities of 2-4k rules. In this paper we present T-Flex (short for flexible TCAM), a technique for implementing virtual TCAM for data center switches.

T-Flex can enhance the rule storage capabilities of a switch by up to 10x, and can forward incoming flows at line rate under most circumstances. Moreover, it provides full OpenFlow compatibility for use in SDN environments, and has been tested to work on a data center production grade Intel switch. Designed as an extension to virtual switches, T-Flex can provide the illusion of an unlimited TCAM in a data center switch.

Keywords—Data center networks; Cloud Computing; Switch Design; Packet Classification; Software Defined Networking.

I. INTRODUCTION

Data center networks - including the ones using CloudStack and OpenStack, require an increasingly large number of policies for providing network and infrastructure services to individual tenants and cloud applications. Typically, these policies are enforced through rules which serve a wide collection of management tasks ranging from cloud security and traffic engineering to ensuring Quality-of-Service (QoS) and maintaining fairness among tenants. In modern data center networks - especially those utilizing Software Defined Networking (SDN), these rules can easily run into hundreds of thousands or even millions. The placement and storage of these rules is thus a critical factor in the efficient operation of the DCN. While there have been attempts to find solutions to the placement problem [1], there has been little effort towards ensuring that modern data center switches are indeed capable of storing the required number of rules while still being able to perform their primary functions correctly and efficiently.

Most modern Top-of-Rack (ToR) switches which perform forwarding and policing functions in the data center network utilize a specialized content addressable memory for storing rules. This memory is housed inside the switch ASIC (Application Specific Integrated Circuit), which allows for hardware based forwarding of data packets at line-rate; usually

in the order of 100 Gbps or higher. However, this memory is expensive to manufacture - and hence is limited in size. For instance, top grade switches from most manufacturers include a TCAM (Ternary Content Addressable Memory) which can store a maximum of 2-4K rules. This is quite insufficient for modern DCNs, as tenant applications and data center administrators require a number of rules to be installed at the switch in order to perform a variety of functions. These rules may relate to firewalls and Access Control Lists (ACLs) for enforcing security policies, or they may be forwarding rules for routing policies or other traffic engineering and monitoring uses. Furthermore, in SDN installations like OpenFlow the switch is required to frequently fetch rules for unknown flows from the SDN controller. This involves a 2-way switch-controller communication which not only introduces a delay in processing packets, but can also cause network congestion due to the increase in control traffic.

These factors indicate the inadequacy of switch ASICs for modern data center applications. However, the existence of powerful server-grade CPUs inside these switches provides a way to overcome these limitations without using any additional hardware. In this paper, we introduce the concept of T-Flex; a Flexible TCAM for data center switches. We leverage the presence of on-board CPUs in ToR switches to develop a virtual memory for storing TCAM rules. T-Flex is designed as a software application for data center ToR switches. The major contribution of this software is the virtual extension of the switch TCAM by utilizing the principles of virtual memory as implemented in Operating Systems. Just like the OS uses on-disk virtual memory to extend the amount of memory available to applications, T-Flex offers the illusion of an infinite memory for storing forwarding rules at the switch.

The major contributions of this paper include: (1) we show that switches can be used to store a larger number of rules; effectively increasing the TCAM size, (2) we demonstrate a prototype implementation of a virtual memory mechanism for storing TCAM rules in a convenient format, (3) we show integration with OpenvSwitch [2] which provides support for OpenFlow on an Intel Seacliff switch, and (4) we demonstrate an increase in the switch's rule storage capabilities by a factor of at least 10 without any significant drop in forwarding performance. The rest of the paper is organized as follows: Section II discusses related work. This is followed by design goals and implementation of the T-Flex software in sections III

and IV respectively. We evaluate performance of T-Flex in section V. Section VI concludes the paper and discusses future work.

II. RELATED WORK

The idea of virtualizing switch TCAM was first proposed by Bhattacharya et al. [3]. Here, the authors suggest exploiting the notion of content locality instead of spatial locality, which is used by modern operating systems in virtual memory implementations. However, the authors have not provided a reference implementation using existing data center infrastructure which is what T-Flex sets out to do. Another related paper [1] proposed an automatic rule management and placement system called Virtualized Cloud Rule Information Base (vCRIB) for data centers. vCRIB circumvents the issue of limited TCAM capacity by dynamically placing rules on aggregation/ToR/edge switches.

A parallel development in the area of networking has been the introduction of extremely powerful switching devices for data center networks. Many ToR switches now include server grade CPUs along with a high speed PCIe interface for packet transfer between the ASIC and CPU, allowing for the development of advanced data plane applications that can run on these switches. ServerSwitch [4] for instance, utilizes the switch CPU for implementing a large forwarding table for flow-based forwarding and a deep packet buffer for bursty traffic.

Recent years have seen a great increase in the popularity of Software Defined Networking. SDN aims to introduce network architectures in which the control and data planes are logically separated and network intelligence and state are centralized. This allows a level of programmability and automation in designing networks which was not previously possible. The most widely adopted SDN standard is OpenFlow [5], which introduces a protocol for controller-switch communication and defines the interface between the control and data planes. However, there has been recent interest in exploring and fixing its limitations [6]. These attempts suggest that SDN still has many architectural issues that need to be figured out before it would become ready for widespread adoption.

In addition, there have been few approaches which focus on keeping incoming packets in the fast path i.e. data plane. This allows maximum number of packets to be processed in the TCAM (which typically has single clock-cycle lookups) and a minimum number of packets are sent to the controller for processing. One approach in this regard is DIFANE [7], which directs the packets to intermediate switches that contain necessary flow rules in their respective TCAMs. These rules are spread over the intermediate switches by the use of rule partitioning. T-Flex also uses the rule partitioning technique but keeps the flows in virtual TCAM.

A substantial amount of research has also been carried out in the area of software packet classification algorithms, most significantly HiCuts [8] and Hypercuts [9]. These algorithms provide a way to implement the task of packet classification on a general purpose processor rather than expensive specialized

hardware. While the processing rate of these algorithms still does not match that of TCAMs present in most switches, they are certainly able to perform well enough for most cases.

III. DESIGN GOALS

In building a virtual memory for switch TCAM, a number of considerations need to be made for the design of the system. Compatibility with existing systems, ease of use and practical means of installation are all factors that have been dealt with in the design of T-Flex. The major design goals and considerations are as follows:

A. Hardware

We chose an Intel Seacliff FM6000 series which includes an Intel Xeon multicore CPU and is capable of running a 64-bit Linux operating system. However, any other data center switch which has similar hardware and can run a Linux OS will support T-Flex. Recent years have seen switch manufacturers introduce devices which include powerful CPUs, which can be used to implement advanced control or management layer applications on the switch. T-Flex is designed with this principle in mind, as it makes use of available hardware that will soon be commonplace in data centers.

All the communication between OpenvSwitch and switch hardware is carried out by using the Intel Ethernet switching platform APIs. These include functions for managing the switch control plane by defining forwarding policies, as well as methods for receiving and sending packets via the switch's front panel ports.

B. Software

T-Flex is designed as an extended and fully ported version of OpenvSwitch, including the *netdev* and *ofproto* drivers [10]. This offers full OpenFlow support, and rules can be installed at the switch by the controller and any other applications using the OpenFlow protocol.

The choice of OpenvSwitch was made due to the fact that it is the most popular virtual switch that supports the OpenFlow protocol. It can be easily ported to hardware switching platforms, or used on its own with a hypervisor. The usage of T-Flex is exactly the same as that of OpenvSwitch. The network administrator can install rules manually through its CLI while the controller can interact with the switch through standard OpenFlow procedures. The only difference is that T-Flex inserts rules first into the virtual memory, from where they are brought into the physical switch TCAM according to the currently active flows.

C. Packet Classification

T-Flex has to perform packet classification in software, for packets which do not have a matching rule in the TCAM. This task requires a fast and efficient classification algorithm, and we have used a variant of Hypercuts [9], which allows us to search the entire virtual TCAM quickly and retrieve the complete set of matching rules.

Hypercuts works by creating a tree data structure to store the entire set of rules. The rules are partitioned in the hyperspace and placed on each node of the tree, from where they can be retrieved when needed. A detailed description of how the algorithm works is beyond the scope of this paper, so we can only summarize the performance improvement that this algorithm brings to packet classification. The standard OpenvSwitch classifier uses a linear hash table that has a worst case time complexity of $O(n)$, where n is the number of rules. Hypercuts on the other hand, can find matching rules in logarithmic time, making it much faster and so better suited for use in T-Flex.

IV. IMPLEMENTATION

The core of the T-Flex application is the data structure for storing TCAM rules in software - the virtual TCAM table. It operates inside the T-Flex application, and provides a convenient way to find matching rule(s) and insert or evict rules from the actual TCAM. The operation of the virtual TCAM table, and how it contributes to the working of T-Flex, is explained in the remainder of this section.

A. Virtual TCAM Table

The customized OpenvSwitch forms the central core of T-Flex. This module is responsible for these major tasks:

1. Receive packets from switch ASIC which do not have a matching rule in the TCAM,
2. Perform packet classification on received packets in order to determine appropriate actions for them,
3. Send packets back out, i.e. perform the actions from Step 2,
4. Provide an OpenFlow interface for flow table manipulation.

The second of these tasks is done using the Virtual TCAM Table (VTT), a packet classifier based on the Hypercuts [9] algorithm. In the VTT, rules are stored in leaf nodes of a tree data structure. Each leaf node stores up to 64 rules, although this number can vary for different rule sets. When a VTT search operation receives a packet for which a matching rule needs to be determined, it returns the leaf node which contains the matching rule. In case there is more than one matching rule, the one with higher priority is selected. The output of the VTT search operation, is a set of rules that need to be inserted into the TCAM.

The main principle behind this approach is called “locality of content”, which makes it possible to efficiently maintain a working set (group of rules) in the switch TCAM. This ideal scenario is depicted in Figure 1a, where incoming packets match rules in the switch TCAM. Step 1 in Figure 1a represents an incoming packet which is looked-up in TCAM at step 2 and packet is forwarded to destination port after TCAM hit in step 3. On the other hand, Figure 1b shows a scenario where an incoming packet encounters a ‘TCAM-miss’ i.e. standard OpenFlow operation (labeled 2), which is then explicitly redirected to the switch-CPU (labeled 3). Once in the Switch-CPU a lookup operation is performed to find the matching rule. The rule is then inserted (labeled 4) into the

switch TCAM and packet is forwarded to destination port (labeled 5).

The virtual TCAM increases the number of rules that can be stored at the switch. Currently, the Intel switch we have used for prototype development can store a maximum of 4k ACL rules. The virtual TCAM on the other hand can store up to 40k rules, which is 10 times more than the physical capacity. This improvement is brought about by saving the complete set of rules in software, and only bringing the relevant rules into the physical TCAM as and when required.

B. Packet Walk

Figure 2 shows the packet walk for the virtual extension of the TCAM. The details of each step are as follows:

- 1) Figure 2 – 200: The packet comes in from Host 1.

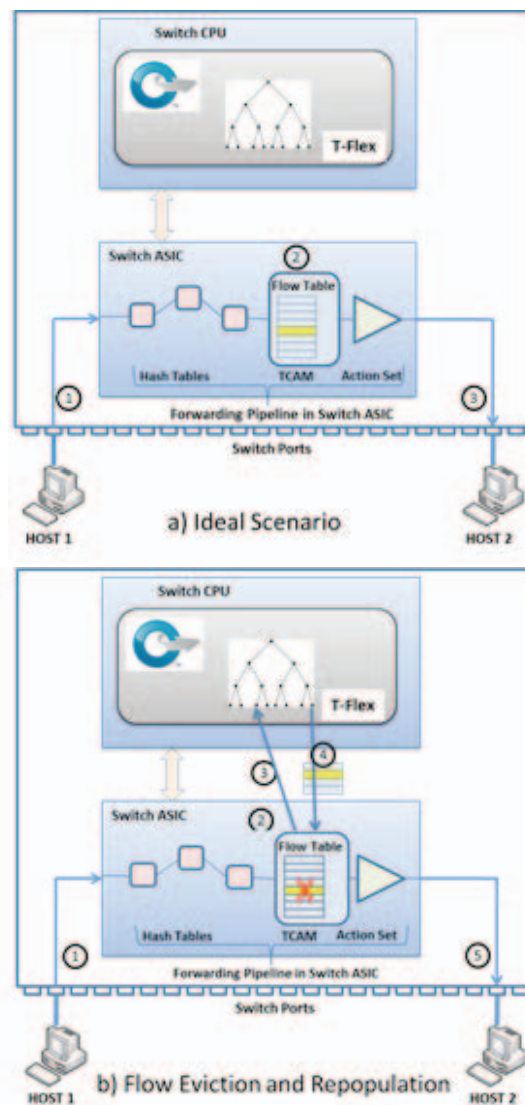


Figure 1: Virtual TCAM operation

2) Figure 2 – 201: The incoming packet has no forwarding matching rule in the switch TCAM and hence a TCAM Miss takes place.

3) Figure 2 – 202: The switch TCAM has a default (lowest priority) rule that redirects all incoming packets via the management link to T-Flex. This particular rule is sticky (stays in TCAM forever) and is installed before executing T-Flex on the switch CPU. All packets with no matching rules in the switch TCAM will hit this rule and get redirected to T-Flex.

4) Figure 2 – 203: A lookup operation to find the matching rule in the VTT takes place. There are two possibilities:

a) Firstly, it may happen that there is no matching rule in the VTT for this incoming packet. In this case, either the default rule is applied or the packet is sent to OpenFlow controller.

b) Second possibility is that a leaf node containing matching rule(s) is found and it does not exist in the switch TCAM. In this case event Figure 2 – 204 takes place and the leaf node is scheduled to be inserted into the switch TCAM.

5) Figure 2 – 205: The Mirror TCAM Table (MTT) is a linked list that mirrors contents of the switch TCAM. Whenever a new node is inserted into the switch TCAM, it is also inserted into MTT. Conversely whenever rules contained in a tree node are deleted from the switch TCAM, they are also deleted from the MTT.

6) Figure 2 – 206: T-Flex application soft switches the incoming packet to its destination front panel port. In this case, since the packet is destined for Host 2, it is sent out from port 2.

7) Figure 2 – 207: T-Flex starts a synchronization thread, which is responsible for inserting all scheduled leaf nodes. While this thread inserts new nodes, it also makes room in the switch TCAM by running an eviction algorithm that inspects corresponding bytes count of each rule and makes the eviction decision. In summary, the sync thread decides to evict any node with least bytes counts since this has been the least used set of rules in the switch TCAM.

C. Partition-with-Splitting

The rules stored in the virtual TCAM usually overlap in one or more fields, and so to ensure correct behavior it is necessary that all over-lapping rules are present in physical TCAM. To summarize, T-Flex must implement the following invariant to work correctly:

Invariant 1: *An incoming packet must match the same highest priority rule when T-Flex is running on the switch CPU with a subset of all rules in the physical TCAM, as compared to when T-Flex is not used and all rules fit in the physical TCAM.*

To implement this invariant, T-Flex relies on a technique called partition-with-splitting [1], which splits leaf node of VTT into independent partitions that can be inserted and evicted to the physical TCAM without any dependencies. This sub-section presents an example to motivate the need for implementing this technique.

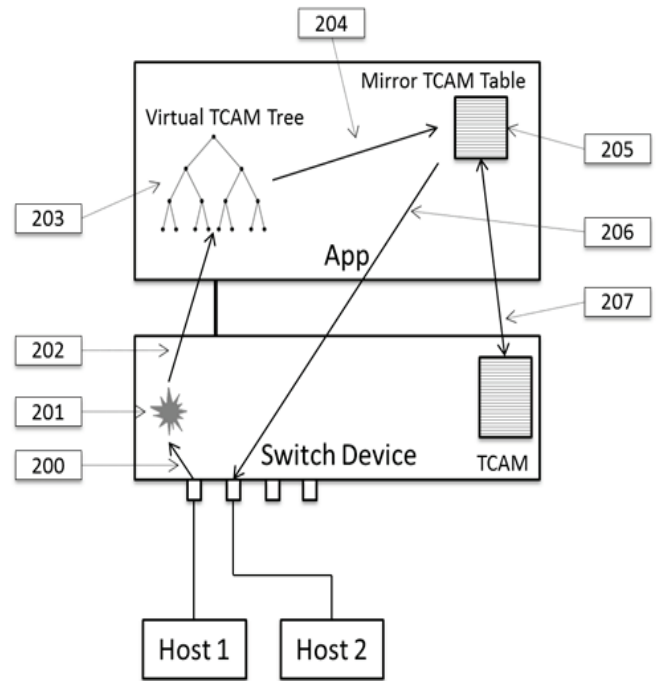


Figure 2: T-Flex packet walk for virtual extension of TCAM

Figure 3 shows an example to motivate the need for implementing the partition-with-splitting technique. For simplicity, we show each rule having two fields—source and destination IP—followed by an action field. Ranges of source and destination IP vary from 0 to 9 while there are only two supported action values ALLOW and DROP. Note that rules with lower index have higher priority i.e. R1 has higher priority than R2. In figure 3 (a), all of the six rules are shown in a 2-D Cartesian grid with source and destination IPs plotted in the x and y axes.

The Hypercuts algorithm builds a tree by cutting the rule space in four equal chunks labeled regions 1 to 4. Figure 3 (a) depicts a cut in the x-axis (source IP) and Figure 3 (b) depicts a cut in the y-axis (destination IP). The resulting tree is shown in Figure 3 (c) where the root node represents the full region and corresponding leaf-nodes representing one of the four sub-regions.

Note that Figure 3(c) has two rules in multiple leaf nodes. For example, rule R2 is contained in leaf nodes corresponding to region #3 and #4. Similarly, rule 6 is contained in leaf nodes corresponding to region #1 and #4. Now imagine a packet comes in the header (1, 2), where 1 and 2 are source and destination IPs. With all four sub-regions in the switch TCAM, this packet matches rule R4 and the packet is thus dropped—because rule R4’s action is DROP. But imagine a scenario where leaf node corresponding to region #1 is evicted out of the switch TCAM. In this case, the same incoming packet matches rule R6 instead and the packet is allowed to switch instead of being dropped—because contrary to rule R4, R6’s action is ALLOW. This is incorrect behavior, which violates invariant 1 stated above.

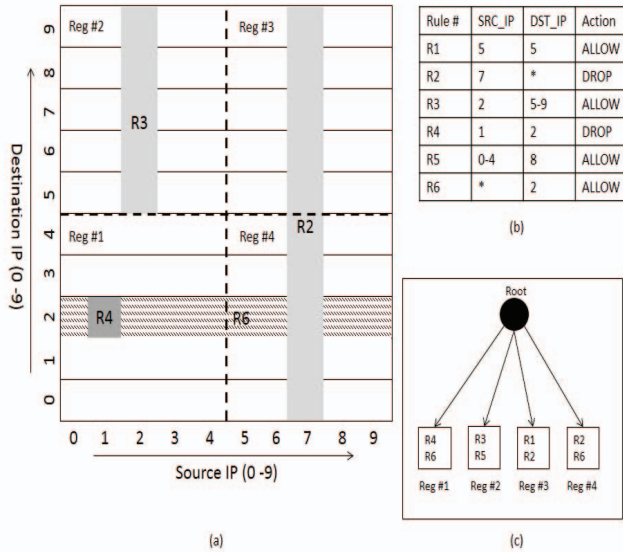


Figure 3: Rule partitioning.

Figure 4 solves the incorrect behavior by using the partition-with-splitting approach. Figure 4 (a) shows list of rules and Figure 4 (b) shows the corresponding tree. On the other hand, Figure 4 (c) and Figure 4 (d) show list of partitioned rules and corresponding tree, respectively. Note that rules R2 and R6 belong to different leaf nodes in the tree. For this reason, both of these rules are split into two rules. For example, rule R2 is split into R2-1 and R2-2 where the former now belongs to region #4 and the latter is part of region #3. Similarly rule R6 is now split into R6-1 and R6-2.

Now imagine if the same packet with header (1, 2) arrives. In case leaf node corresponding to region #1 is present in switch TCAM, then the incoming packet matches rule R4. In case this leaf node has been evicted, then the incoming packet does not match any rule and leaf node for region #1 must be inserted as a result of processing this packet header.

V. EVALUATION

To evaluate the performance of T-Flex, we deployed it in an actual data center environment. The complete setup consists of a T-Flex switch; connected to servers which are used for sending and receiving traffic. For simplicity, we can imagine the same scenario as the one shown in Figure 2, but with end-hosts connected to multiple ports of the switch. The Intel Seacliff Switch provides front panel ports that are capable of switching up to 10 Gbps (Gigabit per second). To take advantage of this, we equipped the hosts with hardware NICs that can transmit and receive traffic at 10 Gbps as well. Packet generation tools Ostinato and pktgen are used to generate traffic which is sent to the switch, where T-Flex operates as an application running on the switch CPU.

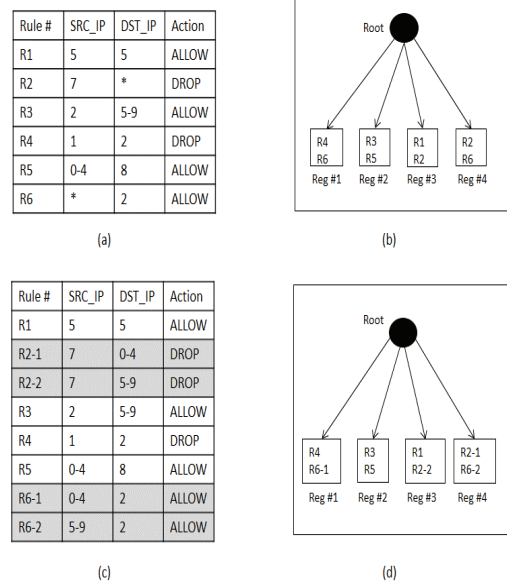


Figure 4: Partition with Splitting

A number of different traffic patterns between these hosts are used to simulate realistic data center traffic scenarios. The evaluation parameters are basically concerned with how the switch performs when a large number of new flows are arriving at its input ports. New flows are defined as those flows which do not have a corresponding rule in the switch TCAM; and so the switch needs to perform a lookup in the virtual TCAM in order to determine the appropriate action for them. Our measurements relate to how varying number of rules affects the forwarding performance of the switch.

To evaluate the performance of T-Flex, we looked at the absolute maximum number of new flows that the switch can forward without packet loss, over a sustained period of time. Next, we evaluated the correctness of operation with VTT enabled, and finally we looked at the performance of the T-Flex switch on real data center traffic traces. The evaluation can be summarized in terms of the following results:

Firstly, T-Flex enhances the switch's capability to handle new flows. Prior work has shown that data center traffic is bursty in nature and follows heavy tailed distribution that can be predictable at the short scales of 1-5 seconds [11]. There is a need of developing fine-grained techniques that can leverage this short-term variability of traffic. The T-Flex switch is able to forward a maximum of 12,000 new flows per second for a duration of 10 seconds without packet loss. This is more than enough to cope with heavy spikes in traffic that are experienced in current data centers. For flow rates greater than this, packet loss is observed due to the fact that there is a large number of ongoing virtual TCAM search operations which increase the processing time for an incoming packet.

Secondly, T-Flex brings about an effective increase in TCAM size. We have shown an increase in effective TCAM size by a factor of 10, and demonstrated the correct operation of the switch with upto 40k rules in the virtual TCAM. This tenfold increase in TCAM size is a function of the memory available to the switch CPU; for switches with greater memory the size of virtual TCAM can be increased further. All flows are forwarded/blocked as per defined rules and no packet loss is observed for upto 12,000 unique flows per second. Even after a 10 folds increase in the virtual TCAM space, there have been negligible increase in the latency of rules installation in switch TCAM. The observed behavior of rules installation in virtual TCAM is exactly same as switch TCAM without using T-Flex.

Finally, T-Flex is fully capable of handling real data center traffic. Our measured performance with traces captured from real data centers [12] were sent to the T-Flex switch and the it was observed that the correct decisions were made for all the flows while minimal packet drops occurred. This translates as the switch being perfectly able to handle all the traffic while maintaining the virtual TCAM which provides a much larger storage for rules.

The challenge in evaluating the performance of a system like T-Flex is that in order to get a real measure of its benefits, it should be deployed in an actual data center setting. As this was not possible for us at this stage, we carried out our evaluation using standard lab tools. This approach has certain limitations in terms of measuring latency and not being able to generate patterns of traffic that are seen in actual data centers. Therefore, we intend to take the evaluation of T-Flex forward by carrying out deployment in a physical data center, where these limitations can be overcome.

VI. CONCLUSION

As demonstrated, T-Flex enhances the capabilities of a data center switch by introducing a “flexible” TCAM. The T-Flex switch works just like a normal data center switch; but can cater for scenarios where a large number of active flows are present which require a proportionally large number of rules. T-Flex has been designed to ensure compatibility with existing data center architectures and is closely aligned with the concept of Software Defined Networking. We have demonstrated the working prototype on a production grade ToR switch, and the modular design of the software also means that it can be easily ported onto other platforms.

We hope that approaches like T-Flex will, in the future be used for a variety of data plane processing applications. Currently, we are looking at more extensive evaluation of T-Flex and utilizing T-Flex for elephant/mice flow classification; in which only elephant flows have forwarding rules in the physical TCAM while mice flows are classified entirely in software. Other areas of interest include alternative architectures where, due to a larger memory of forwarding rules, the switch does not need frequent communication from SDN controllers.

REFERENCES

- [1] Masoud Moshref, Minlan Yu, Abhishek Sharma, and Ramesh Govindan. 2013. “Scalable rule management for data centers,” in proceedings of the 10th USENIX conference on Networked Systems Design and Implementation (NSDI '13), Nick Feamster and Jeff Mogul (Eds.). USENIX Association, Berkeley, CA, USA, 157-170.
- [2] Open vSwitch website, <http://www.openvswitch.org>.
- [3] Suparna Bhattacharya and K. Gopinath. 2011. “Virtually cool ternary content addressable memory,” in proceedings of the 13th USENIX conference on Hot topics in operating systems (HotOS'13). USENIX Association, Berkeley, CA, USA, 3-3.
- [4] Guohan Lu, Chuanxiong Guo, Yulong Li, Zhiqiang Zhou, Tong Yuan, Haitao Wu, Yongqiang Xiong, Rui Gao, and Yongguang Zhang. 2011. “ServerSwitch: a programmable and high performance platform for data center networks,” in proceedings of the 8th USENIX conference on Networked systems design and implementation (NSDI'11). USENIX Association, Berkeley, CA, USA, 15-28.
- [5] Nick McKeown, Tom Anderson, Hari Balakrishnan, Guru Parulkar, Larry Peterson, Jennifer Rexford, Scott Shenker, and Jonathan Turner. 2008. “OpenFlow: enabling innovation in campus networks,” SIGCOMM Comput. Commun. Rev. 38, 2 (March 2008), 69-74.
- [6] Jeffrey C. Mogul, Jean Tourrilhes, Praveen Yalagandula, Puneet Sharma, Andrew R. Curtis, and Sujata Banerjee. 2010. “DevoFlow: cost-effective flow management for high performance enterprise networks,” in proceedings of the 9th ACM SIGCOMM Workshop on Hot Topics in Networks (Hotnets-IX). ACM, New York, NY, USA, , Article 1 , 6 pages.
- [7] Yu, M., Rexford, J., Freedman, M. J., & Wang, J. 2010. “Scalable flow-based networking with DIFANE”. *ACM SIGCOMM Computer Communication Review*, 40(4), 351-362.
- [8] Pankaj Gupta and Nick McKeown. “Packet classification using hierarchical intelligent cuttings,” in proceedings of 7th IEEE Symposium on High Performance Interconnects, 1999.
- [9] J. Sumeet Singh, Florin Baboescu, George Varghese, and Jia Wang. 2003. “Packet classification using multidimensional cutting,” in proceedings of the 2003 conference on Applications, technologies, architectures, and protocols for computer communications (SIGCOMM '03). ACM, New York, NY, USA, 213-224.
- [10] “How to Port Open vSwitch to New Software or Hardware”, available online at: <http://git.openvswitch.org/cgi-bin/gitweb.cgi?p=openvswitch;a=blob;f=PORTING>.
- [11] K. Park and T. Tuan. 2000. “Performance evaluation of multiple time scale tcp under self-similar traffic conditions”. In *ACM Trans. Model. Comput. Simul.*
- [12] Eason, BTheophilus Benson, Aditya Akella, and David A. Maltz. 2010. “Network traffic characteristics of data centers in the wild,” in proceedings of the 10th ACM SIGCOMM conference on Internet measurement (IMC '10). ACM, New York, NY, USA, 267-280.