

Towards Elastic and Resilient In-Network Computing

Daehyeok Kim

CMU-CS-21-143

November 2021

Computer Science Department
School of Computer Science
Carnegie Mellon University
Pittsburgh, PA 15213

Thesis Committee:

Srinivasan Seshan, Co-chair

Vyas Sekar, Co-chair

Justine Sherry, Carnegie Mellon University

Jennifer Rexford, Princeton University

Jitendra Padhye, Microsoft

*Submitted in partial fulfillment of the requirements
for the degree of Doctor of Philosophy.*

Copyright © 2021 **Daehyeok Kim**

This research is sponsored by National Science Foundation grants CNS-1700521 and CNS-1513764, by the CONIX Research Center, one of six centers in JUMP, a Semiconductor Research Corporation (SRC) program sponsored by DARPA, and by a generous contribution from the Microsoft Research PhD Fellowship program.

The views and conclusions contained in this document are those of the author and should not be interpreted as representing the official policies, either expressed or implied, of the United States Government or any other supporting entity.

Keywords: In-network computing, Programmable networks, Programmable networking hardware, Programmable data planes, Resource elasticity, Fault resilience

*To my parents,
for never putting any limits on my potential!*

Abstract

Recent advances in programmable networking hardware technology such as programmable switches and smart network interface cards create a new computing paradigm called in-network computing. This new paradigm allows functionality that has been served by servers or proprietary hardware devices, ranging from network middleboxes to components of distributed systems, to now be performed in the network. The demand for higher performance and the commercial availability of programmable hardware have driven the popularity of in-network computing.

While many recent efforts have demonstrated the performance benefit of in-network computing, we observe a significant gap between what it offers today and evolving application demands. In particular, we argue that in-network computing lacks resource elasticity and fault resiliency which are essential building blocks for practical computing platforms, limiting its potential. Elasticity can address the shortcoming that today's in-network computing only supports a simple deployment model where a single application runs on a single device equipped with fixed and limited resources. Similarly, fault resiliency is critical for managing prevalent device failures for the correctness and performance of applications, but it has gained little attention. Although resource elasticity and fault resiliency have been extensively studied for traditional CPU server-based computing, we find that enabling them on programmable networking devices is challenging, especially due to their low-level abstractions, hardware constraints, heterogeneity, and workload characteristics.

In this thesis, we argue that by designing high-level abstractions and runtime environments that help leverage compute and memory resources available outside of one type of device, we can make in-network computing more elastic and resilient without any hardware modifications. This concept, which we call device resource augmentation, is a key enabler for resource elasticity and fault resiliency for stateful in-network applications written for programmable switches. In particular, we design three systems, named TEA, ExoPlane, and RedPlane, that use this concept to support elastic memory and elastic compute/memory, and fault resiliency, respectively. Each of these systems consists of a key abstraction, programming APIs, and a runtime environment. We demonstrate their feasibility and effectiveness with prototype implementations and evaluations using various in-network applications. Putting all the pieces together, developers can easily enable resource elasticity and fault resiliency for their applications without worrying about underlying complexities.

Acknowledgments

First and foremost, I would like to express my gratitude to my advisors, Professors Vyas Sekar and Srinivasan Seshan. Having Vyas and Srini as my co-advisor is one of the best decisions I have ever made. They gave me the right amount of freedom and guidance during my graduate studies. I am especially thankful for their honest feedback that improved me as a researcher and a person. I am privileged to work with both of them and benefit from both their vision and wisdom. I hope and look forward to continued collaboration with them in the future.

I would also like to thank the rest of my thesis committee members: Professor Justine Sherry, for always being accessible and willing to offer her time to discuss both technical and non-technical topics; Professor Jennifer Rexford, for inviting me to present my work at Princeton University and providing me with detailed and insightful feedback on an early draft of this document; Dr. Jitendra Padhye, for mentoring me since my first internship at Microsoft Research in 2012 and introducing and encouraging me to tackle challenging real-world problems. I am honored to have esteemed and renowned researchers on my thesis committee. In addition, even though they are not officially on my committee, Professor James Hoe and Dr. Victor Bahl have made significant contributions to my thesis and to my overall development as a researcher. I appreciate their willingness to listen to my ideas, support, and make them happen.

Throughout my career, I have had the pleasure to work with several researchers and faculty members: Sharad Agarwal, Anirudh Badam, Changhoon Kim, Jeongkeun Lee, Hongqiang Liu, Jacob Nelson, Sungwoo Park, Dan Ports, Lili Qiu, Shachar Raindel, Insik Shin, Kun Tan, Rashmi Vinayak, and Yibo Zhu. I truly enjoyed our collaborations and appreciated their support and guidance at different stages in my career.

I would like to thank the research and support staff in the Computer Science Department and CyLab: Deb Cavlovich, Catherine Copetas, Chad Dougherty, Diana Hyde, Chelsea Mendenhall, Kathy McNiff, Jamie Scanlon, and AnnMarie Zanger. I appreciate their efforts in making sure everything just works and minimizing the pain of getting a PhD.

I have been fortunate to meet and interact with several talented and fun fellow graduate students and friends: Anup Agarwal, Byeongjoo Ahn, Nirav Atre, Christopher Canel, Shuoshuo Chen, Yong-Seok Heo, Haewon Jeong, Anuj Kalia, Aqsa Kashaf, Dohyun Kim, Juyong Kim, Jack Kosaian, Nikita Lazarev, Jay-Yoon Lee, Robert MacDavid, Kiwan Maeng, Antonis Manousis, Matthew Mukerjee, Daye Nam, Hun Namkung, Pratyush Patel, Devdeep Ray, Zinan Lin, Arjun Singhvi, Milind Srivastava, Ranysha Ware, Hongseung Yeon, Yucheng Yin, Minji Yoon, and Zhipeng Zhao. Thank you for making my journey of getting a PhD more enjoyable.

Finally, I would like to express my greatest gratitude to my family. I sincerely appreciate their many years of unconditional love, support, and encouragement.

Contents

1	Introduction	1
1.1	Today's Myopic View of In-Network Computing	2
1.2	A Vision for Elastic and Resilient In-Network Computing	4
1.3	Summary of Contributions	5
1.4	Dissertation Plan	7
2	Background and Related Work	9
2.1	Evolution in the Control Plane	10
2.2	Evolution in the Data Plane	11
2.3	Programmable Data Plane Technology	12
2.4	Background on In-Network Computing and Applications	14
2.4.1	Justification of In-Network Computing	14
2.4.2	A Taxonomy of In-Network Applications on Programmable Switches	15
3	Table Extension Architecture for Memory-Intensive In-Network Applications	17
3.1	Motivation	19
3.2	Overview	20
3.2.1	Design Space	20
3.2.2	Design Challenges	21
3.3	TEA Design	23
3.3.1	DRAM Access in the Data Plane	24
3.3.2	TEA-Table: Lookup Table Structure	25
3.3.3	Multiple DRAM Servers	30
3.3.4	High Availability	31
3.3.5	Putting It All Together	31
3.4	Implementation	32
3.4.1	Data and Control Plane	32
3.4.2	Programming Network Functions with TEA	33
3.4.3	Limitations	35
3.5	Evaluation	36
3.5.1	Microbenchmarks	36
3.5.2	Application Performance	39
3.5.3	TEA ASIC Resource Usage	40
3.6	Discussion	41
3.7	Related Work	42
3.8	Summary	43

4	On-Rack Switch Resource Augmentation to Support Multiple Concurrent In-Network Applications	45
4.1	Motivation	46
4.1.1	Primer on Stateful In-Switch Applications	46
4.1.2	Motivation	47
4.1.3	A Case for On-Rack Switch Resource Augmentation	49
4.2	Overview	49
4.2.1	Choice of Operating Model	50
4.2.2	ExoPlane Architecture	50
4.2.3	Design Challenges	52
4.3	ExoPlane Runtime Environment	52
4.3.1	Packet-pinning Operating Model	52
4.3.2	Handling Workload Changes	55
4.3.3	Synchronizing Shared Stateful Objects	56
4.3.4	Scaling to Multiple Devices	59
4.3.5	Handling Failures	59
4.4	ExoPlane Planner	59
4.4.1	Inputs	60
4.4.2	Profiler	60
4.4.3	Resource Allocation	61
4.4.4	Application Merger	62
4.5	Implementation	63
4.6	Evaluation	63
4.6.1	Performance in Steady State	64
4.6.2	Performance under Dynamic Workload	67
4.6.3	Shared Stateful Object Synchronization	68
4.6.4	Failover	68
4.6.5	Runtime Resource Overheads	69
4.6.6	Performance of the ExoPlane Planner	70
4.7	Related Work	70
4.8	Summary	71
5	Supporting Fault-Tolerance for Stateful In-Network Applications	73
5.1	Motivation	75
5.1.1	Impact of Switch Failures	76
5.1.2	Existing Approaches and Limitations	77
5.2	Overview	79
5.2.1	Design Challenges	81
5.2.2	Key Ideas	81
5.3	Correctness Model	82
5.3.1	Preliminaries	82
5.3.2	Linearizable Mode	83
5.3.3	Per-flow Linearizability	83
5.3.4	Bounded-inconsistency mode	84
5.4	RedPlane Design	84
5.4.1	Basic Design	85
5.4.2	Sequencing and Retransmission	87

5.4.3	Lease-based State Ownership	88
5.4.4	Periodic Snapshot Replication	89
5.4.5	Protocol Correctness	90
5.5	Implementation	90
5.6	Evaluation	93
5.6.1	Latency in Normal Operation	94
5.6.2	Bandwidth Overheads	95
5.6.3	Failover and Recovery	98
5.6.4	RedPlane Switch ASIC Resource Usage	98
5.7	Related Work	99
5.8	Summary	100
6	Conclusions	101
6.1	Putting All The Pieces Together	102
6.2	Lessons Learned	102
6.3	Future Directions	103
A	Simplified Codes of TEA-enabled NF Implementations	107
B	Lazy Snapshotting Algorithm in RedPlane	109
C	TLA+ Specification of RedPlane Protocol	111
	Bibliography	119

List of Figures

1.1	Today's view of in-network computing: A single application instance is deployed on a single device instance.	3
1.2	Contributions of this thesis: Each system consists of a new abstraction, a programming API, and an runtime environment.	6
2.1	Evolution of programmability in the network control and data planes from industry and the research community over the past 25 years.	9
2.2	Conceptual view of RMT-based programmable switch architecture.	13
3.1	Comparison between RPC-based naïve design and TEA to access external DRAM.	22
3.2	NFs implemented in P4 can be extended with TEA P4 API to look up tables across external DRAM and on-chip SRAM. The control plane is (dotted lines) involved when establishing a TEA channel.	23
3.3	Switch ASIC generates RDMA requests by adding RoCE headers on incoming packets and parse RDMA responses without specialized capabilities for RDMA. To maintain reliable channels, the ASIC maintains per-QP and per-server metadata.	24
3.4	Cuckoo hashing and bounded linear probing. In this example, there are 6 buckets and 2 cells per bucket. The numbers on the top and right side indicate cell and bucket indices respectively.	26
3.5	Design of TEA-Table with scratchpads. Scratchpads temporarily store original packets during lookups. i^{th} bucket of the shadow table has a copy of $((i + 1) \bmod n)^{th}$ bucket of the original table ($n = 6$ in this figure).	27
3.6	Summary of key components in TEA. The components form one logical TEA component (dotted-red box) used by an NF pipeline.	32
3.7	A template of P4 program using TEA abstraction. TEA exposes as a library of P4 control functions (e.g., <code>lookup_response_handler</code>).	34
3.8	Comparison of simplified NAT data plane with and without TEA. To use TEA, in addition to the original logic (white-boxes), developers need to add TEA modules (blue-boxes) and provide basic information necessary for lookup (red-colored).	35
3.9	Lookup performance of TEA via an RDMA channel with a single server.	37
3.10	Scalable lookup throughput of TEA with multiple servers and cache.	38
3.11	Lookup throughput changes during failover events.	38
3.12	Performance of NAT using TEA.	40
4.1	An abstract P4 application and runtime model. An application consists of multiple stateful objects (white boxes) and the control plane logic (blue arrows).	47

4.2	SRAM requirements (normalized to the total amount of SRAM on a switch) with varying workload size and number of applications. If the requirement exceeds 1, it is an infeasible case.	48
4.3	Overview: Green boxes represent the inputs for ExoPlane, and the yellow and blue box indicate key components in ExoPlane.	51
4.4	Inefficient state placement can lead an external device be overloaded, and per-object runtime components incurs high resource overhead. The arrow widths indicate relative flow volumes.	53
4.5	ExoPlane runtime environment processes the majority of traffic at the switch and the remaining at the external device in a run-to-completion manner. The green box is a per-app ExoPlane flow manager, and <i>UKey</i> indicates a union key of the application.	53
4.6	Skewness in flow key (IP 5-tuple): For both Internet backbone and data center case, a few popular keys serve the most of the traffic. This is consistent across measurement epochs.	54
4.7	Basic workflow for inserting new flow entries in ExoPlane.	54
4.8	Incorrect state eviction: application's state has been removed while there is a packet being processed.	55
4.9	Correct two-phase state eviction.	56
4.10	Our state synchronization protocol synchronizes two copies of an entry in the packet counter.	57
4.11	Example inputs for an application p	60
4.12	Merging multiple P4 programs into a single program.	62
4.13	Per-packet processing latency distribution of applications concurrently running on ExoPlane in steady state.	65
4.14	Throughput of each application running on ExoPlane in steady-state with a single external device (Applications are running concurrently).	65
4.15	Scalable throughput with multiple devices. Each color represents a different fraction of traffic offloaded to external devices.	66
4.16	Throughput changes due to workload changes.	67
4.17	Difference in shared object values on the switch and external devices; there are no more packets after the 32 nd epoch.	68
4.18	TCP throughput changes during failover and recovery.	69
4.19	Elapsed time for the resource allocator.	70
5.1	Impact of switch failures on in-switch NATs.	76
5.2	Highlighting why adapting existing approaches for fault tolerance fails for hardware switches.	78
5.3	RedPlane overview highlighting extensions to traditional workflows for in-switch applications	80
5.4	RedPlane state replication protocol packet format.	85
5.5	Basic workflow of RedPlane state replication protocol. "Repl" indicates a state replication request. pkt_n^f indicates n^{th} packet of a flow f	86
5.6	Serializing out-of-order requests with sequencing. Counter values (cnt) in red and blue indicate the state on the switch and the state store, respectively.	87
5.7	Consistent state access for multiple switches.	89
5.8	The main part of P4 implementation of RedPlane-enabled NAT.	92

5.9	Three-layer network testbed for experiments.	93
5.10	End-to-end RTT when RedPlane-NAT processes packets <i>vs.</i> other approaches. . . .	94
5.11	End-to-end RTT for RedPlane-enabled applications. All applications have chain replication enabled for the state store. For Sync-Counter, we also show its overhead without chain replication.	95
5.12	RedPlane replication bandwidth overhead.	96
5.13	Impact of the frequency of snapshotting on RedPlane-enabled HH-detector.	96
5.14	Impact of RedPlane on data plane throughput of RedPlane-enabled applications. .	97
5.15	Impact of update ratio on data plane throughput of the RedPlane-enabled key-value store.	97
5.16	End-to-end throughput changes during failover and recovery with and without RedPlane.	98
5.17	Switch packet buffer occupancy due to request buffering.	98
A.1	Firewall.	107
A.2	Load balancer.	107
A.3	Network address translator.	108
A.4	VPN gateway.	108

List of Tables

- 3.1 Comparison of NF deployment options. We excerpt the information from product briefs [10, 28, 53] and prior work [153, 163]. 19
- 3.2 The NFs we developed with TEA. Table sizes are estimated by assuming 10 million entries with IPv6 addresses. 33
- 3.3 Throughput and latency of NFs implemented using TEA with a single server and corresponding software implementations running on a single server (4 CPU cores). Note that TEA does not involve the CPU on the server. 39
- 3.4 Additional switch ASIC resources used by TEA. 41

- 4.1 P4 applications deployed in a front-end switch of the data center in our motivating scenario. 48
- 4.2 Switch programs written in P4 used in the evaluation in addition to ones introduced in Table 4.1. 64
- 4.3 Comparison of aggregate throughput of four applications running on the appinning model and ExoPlane with four external devices. 67

- 5.1 Examples of stateful in-switch applications and impact of switch failures. 76
- 5.2 Switch ASIC resources used by RedPlane. 99

Chapter 1

Introduction

The traditional view of the network as a “dumb pipe” that merely transmits bits of data from one application end-point to the other is long gone. Today, the network, especially its data plane, becomes more programmable, allowing us to implement sophisticated functionality beyond packet forwarding. The key enabler for this is recent technology advances in programmable network data plane devices. For example, programmable switching Application-Specific Integrated Circuits (ASICs) provide limited data plane programmability while ensuring packet processing rates of 10s of Tbps or a few billion packets per second [21, 50, 53, 54], and programmable Network Interface Cards (NICs), also known as smart NICs, with Network Processing Units (NPUs) or Field Programmable Gate Arrays (FPGAs) support data plane programmability [27, 29, 30], along with the development of target-independent data plane programming languages such as P4 [31, 72] and NPL [43].

This technology advance creates a new computing paradigm called *in-network computing*, which enables us to run various functionalities that traditionally have been served by CPU servers or proprietary hardware devices, ranging from network middleboxes to distributed systems (*e.g.*, [78, 115, 116, 148, 153, 162, 186, 197]) on network data-plane devices. Essentially, this makes the network itself a new kind of computing platform. We observe that this new form of computing is becoming promising today because user demands and technological advances have converged quickly. Today, while network bandwidth is growing rapidly up to several hundreds of Gbps, the speed of CPU performance improvement has decreased with the slowing of Moore’s law. This can cause the low and unpredictable performance of networked applications, such as high tail latency. On the other hand, programmable data plane devices can process data with limited flexibility at a high rate while providing predictable performance, in-network computing can provide opportunities for improving the performance of the applications, as well as reducing operational costs.

Unfortunately, despite many recent efforts that have demonstrated the promise of in-network computing, we observe that there are still many missing pieces to make the programmable data plane a practical and future-proof computing platform, primarily to support evolving application workloads. In particular, as we will describe later in this chapter, while *resource elasticity* and *resilience against device failures* are essential for any practical computing platforms, they have been considered in an ad-hoc manner or gained little attention in today’s in-network computing.

Although the resource elasticity and fault resilience have been extensively studied in the context of traditional server-based computing, we find that it is uniquely challenging to support them for in-network computing, due to static and limited resource and capability within a single device, limited programmability with low-level primitives, and high-performance requirements.

In this thesis, we argue that by designing abstractions that effectively expose resources and capabilities available *outside* of one type of data plane device (*e.g.*, programmable switches) while hiding the complexities of dealing with heterogeneity, we can make in-network computing more elastic and resilient without any hardware modifications. This concept, which we call *device resource augmentation*, is a key enabler for resource elasticity and fault resiliency. This is based on our observation that while each type of device has certain limitations in terms of available resources and capabilities, different types of devices can complement each other. For example, while a programmable switch can process packets at the rate of Tbps with only a few 10s MB of SRAM, a smart NIC can provide a few GB of DRAM at the cost of lower processing capacity.

We particularly focus on enabling resource elasticity and fault resilience for in-network applications running on programmable switches, also called in-switch applications, which have the most tight constraint on resource and capability among programmable data plane devices available today, while leveraging resources on other types of *external* devices (*e.g.*, smart NICs and commodity CPU servers). In particular, we design three systems for in-switch applications, each of which consists of a key abstraction, programming APIs, and a runtime environment: (1) **Table Extension Architecture (TEA)** provides a virtual table abstraction that gives an illusion of large match-action tables by effectively utilizing DRAM available on remote servers (*elastic memory*); (2) **ExoPlane** offers an operating system for an on-rack switch resource augmentation architecture to support multiple concurrent applications with an infinite switch resource abstraction (*elastic compute and memory*); and (3) **RedPlane** provides an illusion of one big fault-tolerant switch for an application deployed on multiple switches in the network (*fault resilience*). Putting all the pieces together, they realize elastic and resilient in-network computing without any hardware modifications.

Thesis statement: *With the right abstractions for resources on heterogeneous programmable data planes and runtime environments, in-network computing can be made more elastic and resilient.*

In the remainder of this chapter, we describe today’s limited view of in-network computing, our vision of elastic and resilient in-network computing, challenges in realizing the vision, and then highlight our key results.

1.1 Today’s Myopic View of In-Network Computing

With the technology advance in programmable data planes, many recent efforts have shown the feasibility and effectiveness of various applications running on programmable data plane devices, including network monitoring and telemetry [78, 146, 186], DDoS defenses [148, 197], key-value store caches [115, 144], network middlebox functions [153, 162, 168], consensus protocols [116,

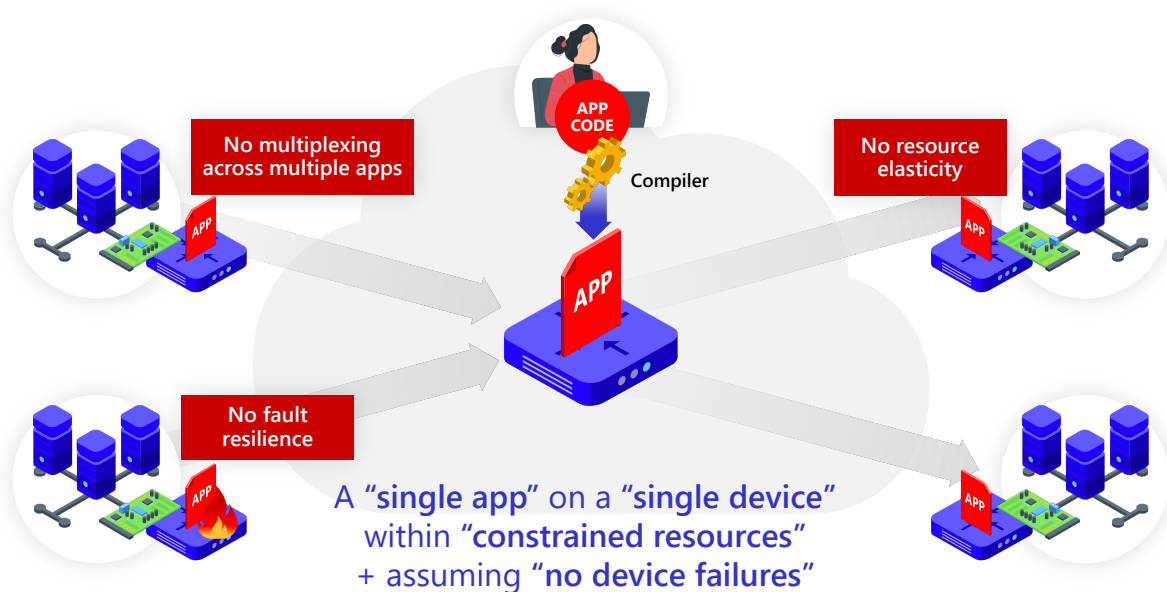


Figure 1.1: Today’s view of in-network computing: A single application instance is deployed on a single device instance.

[137], and others [174, 189]. Today, network administrators or application developers program the data plane devices like in the early days of server-based computing; they program and deploy *a single program on a single device* within a given fixed and limited physical compute and memory resources.

While the current way of enabling in-network applications might be enough for running a single program with fixed and small-sized workloads, we argue that it can significantly limit the potential of in-network computing. First, we observe a huge gap between evolving demands from network administrators and developers and the current *myopic* view of in-network computing, especially due to the lack of *resource elasticity* and *fault resiliency*, as illustrated in Figure 1.1. Although the number of in-network applications that can potentially run concurrently keep increasing [109, 125] and the amount of traffic that needs to be handled also grows [25, 81], today, only a single program can run on a single device under tight resource constraints (*e.g.*, 10s MB of SRAM on a programmable switch [66]) which is typically hard to be extended without significant hardware modifications. Thus, if applications’ workload demands exceed the constrained resource capacity on a single device, the applications will fail to run. Second, in addition to the resource constraint, several measurement studies in production networks have shown the prevalence of networking hardware failures [99, 143, 152], which can affect the performance and correctness of applications [130]. However, handling such failures has received little attention or considered in an ad-hoc manner.

In summary, although resource elasticity and fault resilience are fundamental building blocks for any practical and future-proof computing platform, they have not been considered in today’s in-network computing, significantly limiting its potential.

1.2 A Vision for Elastic and Resilient In-Network Computing

In this thesis, we envision an in-network computing platform that natively supports resource elasticity and fault resilience for in-network applications, rather than letting network administrators and developers constrain their workload sizes (in terms of the number of applications and traffic volumes) or manually handle failures in application-specific ways.

Supporting resource elasticity and fault resiliency for a computing platform is not a new problem; they have been extensively studied for traditional server-based computing, especially in the context of operating systems and distributed systems. For example, while in the early days of server-based computing, people had to write a program within limited memory space, with the advent of virtual memory abstraction offered by operating systems [63], they can program more easily with an illusion of large memory space given the limited physical memory space. Similarly, while handling failures can be very complicated and error-prone, abstractions such as replicated state machines [158, 176] provide an illusion of a single, highly available system. Such abstractions enable elastic and resilient computing while hiding the complexities of how it actually works. So then, the natural question is: *Can we apply these existing techniques designed for server-based computing to in-network computing to achieve resource elasticity and fault resiliency?*

Unfortunately, we find that there are unique practical challenges in supporting the resource elasticity and fault resiliency for in-network computing platforms, which makes it challenging to apply the existing techniques:

- **Challenge 1: Fixed resource pool.** Unlike CPU servers which have a hierarchy of resources (*e.g.*, CPU caches, main memory, and disks for memory), the current programmable data plane devices, especially programmable switches which are the context of this thesis, only have limited resource and capability within a device, which are fixed at a hardware design time. Thus, there seems no way to utilize resources elastically.
- **Challenge 2: Limited programmability.** Second, limited programmability with low-level primitives such as bit-level packet header manipulation in the current data plane programming languages such as P4 [72] make it hard to implement existing complicated mechanisms in the device data plane, which are designed for more capable CPU servers with higher-level primitives.
- **Challenge 3: High-performance requirement.** Lastly, even if we can somehow implement the mechanisms in the device data plane, since the data plane needs to process a stream of packets at a very high rate (*e.g.*, a few tens Tbps in the switch data plane), it is challenging to run the mechanisms without affecting the performance.

Fortunately, we find an opportunity to address the first challenge raised by the fixed and limited resources on a single device, which seems to be the most fundamental obstacle to realizing our vision. We observe that although each type of device has certain limitations in terms of available resources and capabilities, if we look at *outside* a single device instance, there are other types of data plane devices accessible through the network, and different types of devices can be complementary with each other. For example, while a programmable switch can process packets at the rate of Tbps with only a few 10s MB of SRAM, a smart NIC provides lower processing capacity with a few GB of DRAM. Thus, if we can somehow let applications running on one

type of device (*e.g.*, a switch) utilize larger resources on another type of device (*e.g.*, DRAM on a smart NIC), that can be a potential solution to enable elastic resources and potentially provide a way of supporting resilience as well.

Although this approach sounds promising, realizing it is not trivial, especially due to the two remaining challenges: limited programmability and high-performance requirement. Since applications running on a switch now need to access resources on different types of devices (*e.g.*, from switches to NICs and CPU servers) connected through the network using limited programmability, sustaining high performance while not affecting the correctness of applications and while hiding the complexities from programmers is challenging.

Our high-level approach to address these challenges consists of three steps: First, we understand the requirements and characteristics of application workloads (*e.g.*, memory requirements and state access patterns); second, we carefully design runtime environments that allow a switch to utilize resources on other types of devices (*e.g.*, a smart NIC) to enable resource elasticity or fault resilience, which can be implemented under hardware constraints; and lastly, we design abstractions and programming APIs that expose the new capabilities to application developers and network administrators to implement applications that their applications without worrying about the underlying complexities of accessing and managing resources on heterogeneous devices.

1.3 Summary of Contributions

This thesis presents three novel systems that realize our vision of elastic and resilience in-network computing platform, built based on the above approach. In particular, each of them provides a new abstraction and runtime environment that exposes resources on external devices for in-network applications running on the switches to support resource elasticity and fault resiliency:

Table Extension Architecture (TEA) is a system designed to offer *elastic memory* for state-intensive applications. It consists of a programmable switch and multiple servers offering DRAM and provides an illusion of large match-action tables for applications running on programmable switches via a new abstraction called virtual table abstraction. In TEA, we demonstrate:

- The feasibility of table lookups on available DRAM on remote servers entirely in the data plane with low and predictable latency (1.8–2.2 μs) while not involving CPUs on servers.
- The feasibility of providing resource elasticity in terms of table size and access bandwidth that can be scaled linearly by recruiting more servers (*e.g.*, 138 million lookups per second with 8 servers in our testbed).
- The cost and performance benefits of TEA-enabled applications compared to the same applications running on servers.

We present TEA in [Chapter 3](#).

ExoPlane presents an operating system for on-rack switch resource augmentation architecture designed to offer *elastic compute and memory* for multiple concurrent stateful applications. The architecture consists of a programmable switch connected to multiple external data plane devices

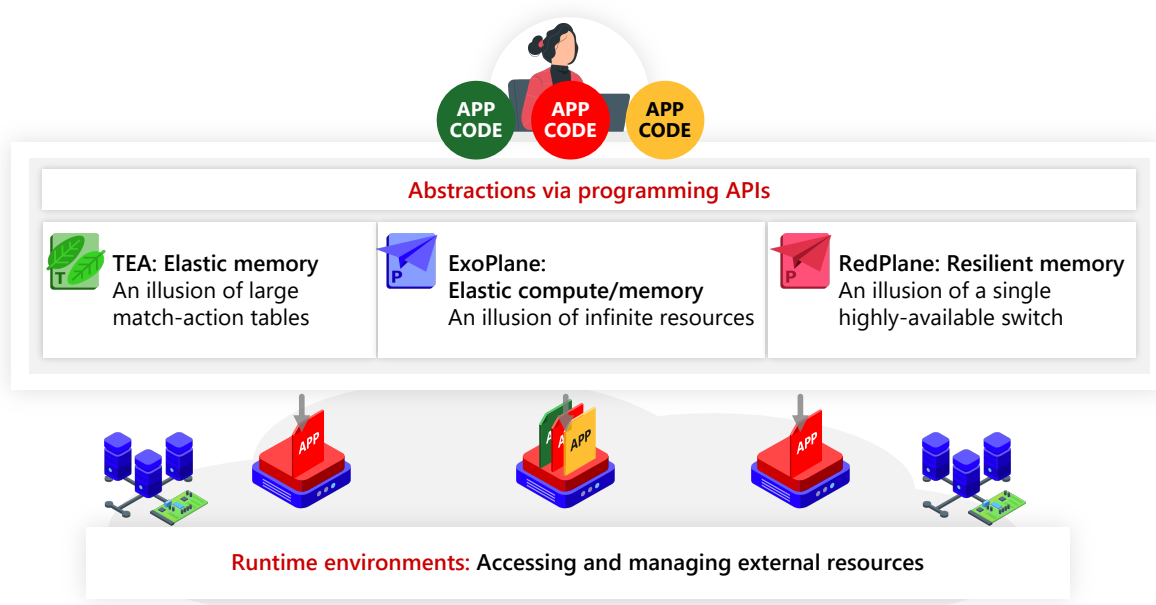


Figure 1.2: Contributions of this thesis: Each system consists of a new abstraction, a programming API, and an runtime environment.

on the same rack. And ExoPlane provides an OS-like abstraction that effectively multiplexes applications across the switch and external devices. In ExoPlane, we show:

- The feasibility of multiplexing concurrent stateful applications written for a switch across the multiple heterogeneous devices (NPU-based smart NICs in our prototype) while providing predictable (*e.g.*, ≈ 300 ns at the switch and 5.5 ns at an external device in steady state) and scalable (*e.g.*, up to 394 Gbps, the maximum rate in our testbed) performance.
- The effectiveness of abstractions for resource and state management that allow to hide the complexities of dealing with heterogeneity from network administrators and application developers.
- The low control and data plane resource overhead while providing the above benefits.

We present ExoPlane in [Chapter 4](#).

RedPlane is a system designed to offer *fault resilience* for applications running on programmable switches. It consists of multiple programmable switches where the same application is running and servers offering in-memory state store. It provides one big fault-tolerant switch abstraction that gives an illusion of a single, highly-available switch. In RedPlane, we present:

- The definition of a new correctness model based on the traditional notion of *linearizability* for in-switch applications and its realization with RedPlane protocol.
- The feasibility of implementing RedPlane protocol entirely in the data plane, which correctly replicating state to external state store for different types of applications (*i.e.*, read-centric, write-centric, and read/write-mixed).

- The performance benefit compared to alternative approaches by showing negligible per-packet latency overhead for read-centric applications and less than $8 \mu s$ overhead even for the worst case and fast recovery of end-to-end TCP throughput within a second.

We present RedPlane in [Chapter 5](#).

Overall, our systems provide key building blocks for enabling elastic and resilient in-network computing, as illustrated in [Figure 1.2](#). Once network administrators or application developers write their switch programs using our APIs (or without any modifications), our runtime environments allow applications to leverage external resources to support large workloads or make them tolerant of switch failures. They enable these without any hardware modifications.

1.4 Dissertation Plan

This thesis proceeds as follows. In [Chapter 2](#), we discuss the background and related work, including the history of programmable networks, programmable data plane technology, and in-network computing and applications. In [Chapter 3](#), we present TEA for state-intensive in-network applications, which provides a virtual table abstraction. In [Chapter 4](#), we discuss ExoPlane, an operating system for an on-rack switch resource augmentation architecture to support multiple concurrent stateful in-network applications. In [Chapter 5](#), we discuss RedPlane, a framework that makes in-network applications tolerant of switch failures by providing a single highly-available switch abstraction. Finally, in [Chapter 6](#), we discuss lessons learned and future research directions and conclude.

Chapter 2

Background and Related Work

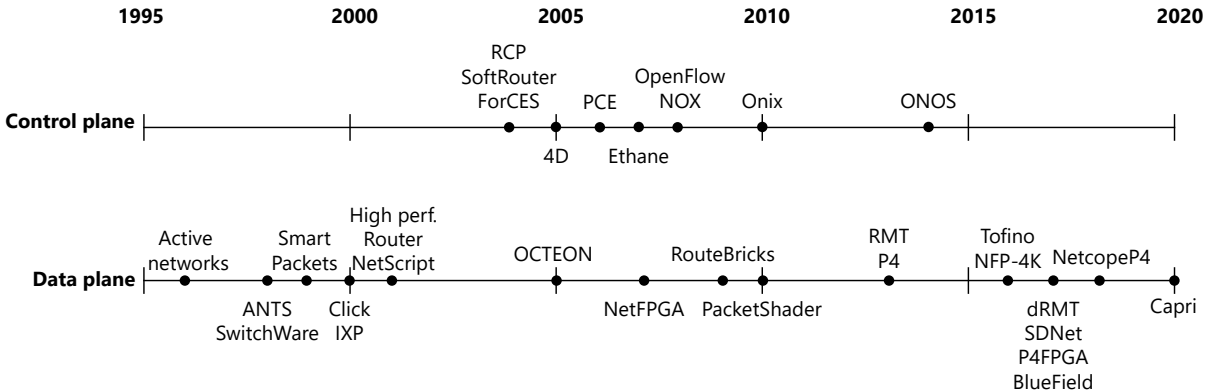


Figure 2.1: Evolution of programmability in the network control and data planes from industry and the research community over the past 25 years.

Over the past two decades, the research community and industry have made various efforts to make computer networks more programmable, primarily to easily manage complex networks and add new functionality. Such efforts toward programmable networks have made indirect or direct impacts on the in-network computing paradigm. In this chapter, we provide the background and related work on programmable networks and in-network computing. In particular, we focus on how the network control and the data plane have been evolved and impacted today’s programmable networks and in-network computing. [Figure 2.1](#) summarizes the development of programmability in the control and data planes over the past 25 years.¹ We see that there have been continuous efforts in both planes simultaneously.

In the following sections, we first discuss the evolution in the control plane ([Section 2.1](#)) and data plane ([Section 2.2](#)), describe details of programmable data plane technology ([Section 2.3](#)), and provide a background on in-network computing and applications ([Section 2.4](#)).

¹This figure is inspired by and partly based on the discussion in Feamster *et al.* [93].

2.1 Evolution in the Control Plane

While it seems natural to have the network control plane as a separate component in today's networks, traditionally, there was a tight integration between the control and data planes. This made it challenging to configure the network, debug network configuration problems, and analyze network-wide routing behaviors, motivating the research community and industry to decouple network devices' control and data planes.

Control-data plane separation (before 2005). There were several early attempts to build logically centralized network controller by moving control functionality to commodity servers and designing interfaces and protocols between the controller and network devices (*i.e.*, switches and routers). It leveraged improved computing power and memory capacity of the servers while providing the controller with a network-wide view. For example, the SoftRouter [133] ran a control software on an external server while allowing it to communicate with the device data plane via the Forwarding and Control Element Separation (ForCES) API [198] to install entries to forwarding tables in the data plane. Although it seemed promising, it required switch and router vendors to modify hardware to support the new API; thus, it has not been widely adopted. As an alternative approach, the Routing Control Platform (RCP) [92] used an existing standard control-plane protocol, the Border Gateway Protocol (BGP) [149], to install forwarding table entries in the data plane of legacy switches and routers without hardware modifications.

Clean-slate architectures (2005 – 2007). While those proposals made some progress in separating the control and data planes, they have not been adopted by equipment vendors and network administrators due to either (1) the requirement of implementing new APIs (*e.g.*, ForCES) or (2) the support for a limited range of applications running over existing routing protocols (*e.g.*, RCP using BGP). It motivated a group of researchers to broaden this vision of separating the control and data planes and design clean-slate architectures for logically centralized controller. The 4D project [101] proposes an architecture consisting of four layers: the decision plane (consisting of logically centralized controllers that convert network-level objectives into packet-processing state), the dissemination plane (for installing packet-processing state to the data plane), the discovery plane (for discovering network elements and creating a logical network map), and the data plane (for handling individual packets based on state given by the decision plane). Following this vision, there was work that proposes a concrete implementation such as the Ethane project [76] and SANE [75], which creates a logically centralized controller deployed in a campus network and uses flow-level access control as an example case.

OpenFlow and network operating systems (2008 – today). Lessons learned from earlier attempts on control-data plane separation and prior efforts on clean-slate architectures and their implementation significantly impacted the OpenFlow [151]-based realization of software-defined networking (SDN). By locating the sweet spot between the vision of high programmability (*e.g.*, managing the network from a logically centralized controller) and pragmatism (*e.g.*, not requiring significant hardware modifications), OpenFlow has gained huge attention and been adopted by industry. Many device vendors had manufactured OpenFlow-compatible switch devices, and many controller platforms that use the OpenFlow API had been developed [38, 103, 132]. Based on the OpenFlow ecosystem, network administrators could start writing a consolidated control

plane program that manages fixed-function switch data planes (*e.g.*, changing packet forwarding rules) and run it on a centralized controller. Each OpenFlow switch has match-action tables, each of which contains a set of rules consisting of match patterns (*e.g.*, match on a specific packet header field) and actions (*e.g.*, drop, flood, or forward a packet, modify packet header fields, and forward the packet to the controller), and the control plane program can install, modify, or delete entries in the tables via the OpenFlow API. While there is one logically centralized controller, multiple physical controllers could run on multiple servers for better scalability, reliability, and performance while providing an illusion of a single controller. To this end, the Onix controller [132] proposes the idea of a network information base that represents the network topology and other control states shared by all controller replicas while supporting state consistency and durability. More recently, the ONOS project [9] offers a distributed controller with similar functionality.

2.2 Evolution in the Data Plane

Despite the evolution in the control plane programmability and the success of the SDN, the data plane programmability remained challenging until recently. This section describes how the data plane has evolved to support programmability, from early efforts on active networks to the recent development of programmable hardware switches and NICs.

Active networks (– early 2000). The initial proposal of the active networks [187] and the subsequent studies [62, 73, 195] were the first attempts to make the network data plane programmable. They proposed programming models and interfaces (*e.g.*, the capsule model [194] and the programmable router/switch model [69, 182]) that exposed compute and memory resources on individual network nodes, which developers can use to implement custom functionality to process a subset of packets passing through the node. The key motivation and enablers behind active networks were (1) the cost reduction in computing that allowed to put more processing in the network, and (2) the technology advances in programming languages such as Java that offered platform portability and safe code execution using virtual machines protecting the host machine and other processes from misbehaving programs [181]. While the idea of active networks sounded promising, it was not adopted by industry and deployed in real networks; the critical stumbling block was the lack of hardware technology for network nodes that can achieve both programmability and high performance simultaneously for various applications, making it difficult to prove the efficiency and effectiveness of the approach beyond the lab environments.

Software routers on commodity servers (2000 – today). At the same time, since the early 2000s, there have been efforts to enable high-performance and modular software router functionality on commodity servers. The Click modular router [131] was one of the early efforts focusing on modularity and programmability. Based on this modular framework, the RouteBrick [87] system proposed an architecture that can achieve high performance and programmability by exploiting parallelism across multiple servers and multiple cores within a server. Moreover, the PacketShader [106] framework further optimized and accelerated the packet processing pipeline by leveraging massive parallelism supported by general-purpose Graphics Processing Units (GPUs).

Programmable hardware switches (2000 – today). In the early 2000s, Intel developed Network Processor Units (NPUs) named IXP [1] to replace the general-purpose embedded microprocessors and ASIC combinations used in network routers. While it provided both performance and programmability somewhere in the middle of the microprocessors and ASICs, its performance was limited to a few 10s of Gbps, which is far below today’s requirement. In 2013, a group of researchers from academia and industry developed a programmable switch ASIC architecture called Reconfigurable Match Tables (RMT) [71], which becomes a basis for the Barefoot’s Tofino chip [53]. With the support for a data-plane programming language such as P4 [72], which was first proposed in 2013, Tofino became one of the widespread programmable switching ASICs used by academia and industry. Around the same time, other vendors, including Intel and Cavium (later acquired by Marvell), released programmable switching ASICs called FlexPipe and Xpliant, respectively [8, 21]. A few years later, another group of researchers developed disaggregated RMT (dRMT) [80] that addressed some architectural restrictions (*e.g.*, tight coupling between memory and stages) of the original RMT architecture.

Programmable hardware NICs (2007 – today). NetFPGA [58] was one of the earlier network-attached FPGA-based platforms that have been used for prototyping network equipment, including the Ethane [76] switch. While the FPGA chip is typically programmed using Verilog or other low-level languages, recently, there has been some work from academia and industry on compiling programs written in high-level languages such as P4 into an RTL to improve programmability [19, 42, 112, 193]. Moreover, Netronome has manufactured NPU-based NICs [30] that can be programmed using P4 and Micro-C. In addition, NVidia Mellanox manufactured another type of smart NIC called Bluefield [34] equipped with multi-core ARM CPUs and their ConnectX ASIC chip. More recently, Pensando systems have released an ASIC-based P4 programmable NIC called Capri [47].

2.3 Programmable Data Plane Technology

We now discuss details of a few representative types of programmable data plane hardware, which are widely used today, including programmable switches and ASIC, NPU, SoC, and FPGA-based smart NICs.

Programmable switching ASICs. Programmable switch architectures used today, *e.g.*, Intel Tofino [53], use a limited amount of on-chip memory (*e.g.*, SRAM and TCAM) to provide a variety of stateful object abstractions, including tables, registers, meters, and counters. Developers can use these abstractions to keep state across multiple packets, such as the address translation table in the network address translators (NATs). In the ingress and egress match-action pipeline, objects are allocated in each stage and accessed by packets via arithmetic logic units (ALUs). These objects are also accessible by the switch control plane through the ASIC-to-CPU PCIe channel, which has a limited bandwidth ($O(10\text{ Gbps})$) compared to the ASIC’s per-port bandwidth ($O(100\text{ Gbps})$). In addition, the ASIC provides other built-in functionality such as packet replication, recirculation, and mirroring for more advanced packet processing.

To understand how a program is mapped to switch pipelines, we focus on a reconfigurable match-action table (RMT) [71]-based programmable switch architecture illustrated in Figure 2.2.

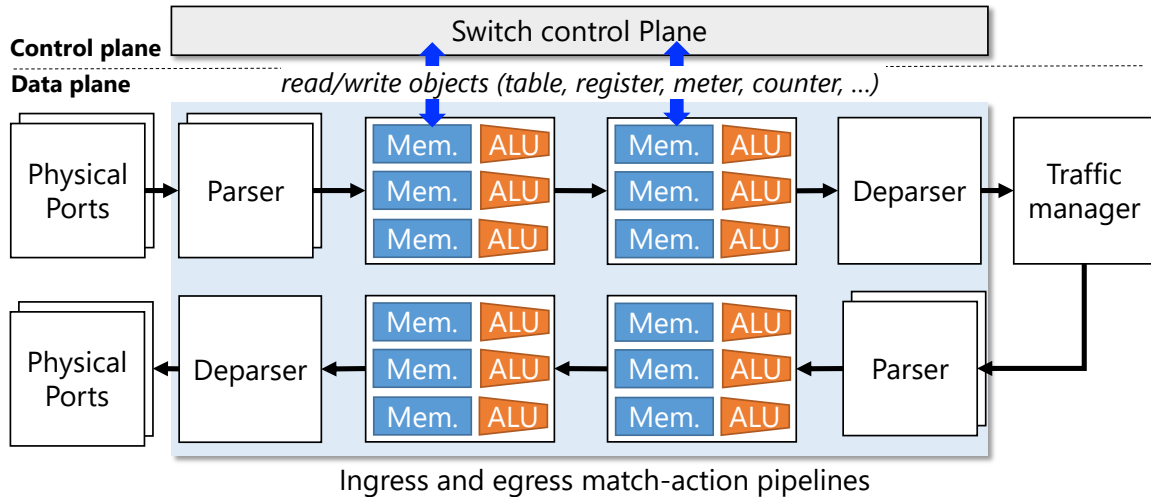


Figure 2.2: Conceptual view of RMT-based programmable switch architecture.

In this architecture, packets are processed in a streaming manner in a hardware *pipeline* consisting of multiple match-action stages with memory and compute resources. When a target-specific compiler (*e.g.*, Tofino P4 compiler [45]) compiles program codes, it places each object to one or more stages based on the resource requirements and the resource availability, and the placement is fixed at compile-time. Typically, the amount of required memory increases proportionally to the size of an object (*e.g.*, the number of table entries or the size of a register array), whereas compute resources are allocated per-object basis (*e.g.*, one stateful ALU is allocated to a register array regardless of its size). If the resource demand from the program exceeds the available resources in the switch, the compilation process will fail.

ASIC-based smart NICs. ASIC-based NICs, such as NVIDIA Mellanox ConnectX-6 [35], provide limited programmability for fixed functions implemented on the ASIC through proprietary APIs. Its pipeline is not fully programmable, but one can add and modify flow rules, similar to configuring OpenFlow-compatible switches. More recently, Pensando systems have released P4-programmable ASIC-based smart NICs called Capri [47], but their performance and programmability are unknown.

NPU-based smart NICs. NPU-based NICs, such as Netronome Agilio CX [30], consist of an array of wimpy packet processors and provide full programmability through languages such as P4 or Micro-C. NPUs are programmable processors optimized for several operations frequently used in packet processing, such as packet I/O, table lookups, queue management, and header manipulation. While NPUs are more programmable than switching ASICs, they do not scale beyond a few 10s of Gbps, which is much lower than switching ASICs (a few Tbps). Most NPU-based NICs support stateful packet processing by providing access to a few 10s of MB of SRAM and a few GB of DRAM for the processors that operate on packets.

Multicore SoC-based smart NICs. Multicore SoC-based NICs (*e.g.*, NVidia Mellanox Bluefield [34]) are equipped with on-board CPU cores alongside the ASIC, which is fully programmable like a regular CPU. These CPU cores are not on the data path and connected with

the ASIC via a PCIe channel, so they are not adequate for processing every data packet due to the performance overhead.

FPGA-based smart NICs. These NICs look like regular FPGA cards but are equipped with Ethernet ports attached to the networks, which is fully programmable either using low-level languages such as Verilog or libraries for high-level languages such as OpenCL or P4 (*e.g.*, Intel N3000 [41], Xilinx U50 [56], and Xilinx U280 [55]). As a result, they can be highly customized to a specific application and potentially achieve speeds as high as 200 Gbps. Many leading cloud service providers such as Microsoft [94] and Amazon [22] have deployed the FPGA-based NICs in their data centers to offload various tasks, including indexing for search engines, machine learning, storage, and virtual network functionality [77, 94].

2.4 Background on In-Network Computing and Applications

The advancement of data plane hardware technology and the development of domain-specific data plane programming languages such as P4 [72] and NPL [43] have become key enablers for in-network computing. Network administrators and developers now can easily implement their custom stateful packet processing logic using the languages and deploy them to the network data plane. In this section, we first discuss when and which applications can most benefit from in-network computing and then describe a taxonomy of existing in-network applications.

2.4.1 Justification of In-Network Computing

Although programmable data plane hardware devices become an attractive platform to serve various applications, we first need to justify whether a given application can benefit from in-network computing, especially compared to running it on end-host CPU servers, before deciding to offload or implement it on the switches or NICs. Based on our analysis of various in-network applications, we observe that considering the following three aspects is necessary for justification:

Functional feasibility. Most importantly, an application or functionality must be implementable on target devices (*e.g.*, switches or NICs). As we described in Section 2.3, each hardware type has different capabilities and resources. For example, programmable switches guarantee the line-rate packet processing speed, but they do not support complex arithmetic operations and deep payload inspection. Given such constraints, network administrators and application developers need to check whether their applications can be realized on target devices.

Operational regimes. Given multiple feasible deployment options (*i.e.*, target devices), network administrators need to choose one by considering the performance, resource constraints, cost, and energy efficiency of each option and their workloads and operating conditions (*i.e.*, traffic rate and the number of concurrent flow that an application instance have to process). For example, while NATs can be implemented on programmable switches or CPU servers, if the traffic volume is small (*e.g.*, a few 10s of Gbps), using a few servers could be a better option than implementing it on a programmable switch, in terms of cost and energy efficiency.

Benefits from “in-network” deployments. Another critical factor to consider is whether an application can benefit from in-fabric or in-network deployments, compare with end hosts or

server-based deployment models. Depending on where the application is deployed, it can have different visibility of the network and performance. For example, running the application at an aggregation layer of a data center network can provide higher visibility than running it on NICs or end-host CPU servers. As we will see later, certain applications, such as network monitoring and security defenses, could benefit from running on programmable switches. Also, when the application runs on programmable switches where all traffic flows through, it incurs little or no additional per-packet latency, compare with running it on NICs or end hosts, which requires re-routing packets and could increase per-packet latency.

2.4.2 A Taxonomy of In-Network Applications on Programmable Switches

Many recent proposals have shown the promise of in-network computing using programmable switches as a platform. Based on our observation above, we classify them into four categories and discuss why they can benefit from in-network computing. In particular, we discuss in-network applications developed for and running on programmable switches, which is the focus of this thesis.

Network middlebox functions. Network middleboxes refer to a device that performs functions other than the standard functions of routing or switching packets between a traffic source and a destination [74]. Examples include NATs, load balancers, firewalls, and intrusion detection systems (IDSes). Recent work has shown the performance and cost benefits of implementing middlebox functions on programmable switches [115, 123, 137, 138, 139, 144, 147, 153, 162, 200]. Switch-based middlebox implementations have advantage over alternative deployment options such as using proprietary hardware devices or CPU-based network function virtualization (NFV), in terms of traffic volume it can support and performance (*e.g.*, a single switch can serve tens of Tbps of traffic at line rate). Also, since switches are typically located on the network path, it does not incur additional latency due to re-routing to middleboxes. However, programmable switches do not support all types of middlebox functions; due to their limited computational capabilities, they are good at serving lookup-heavy, compute-light middlebox functions. They cannot effectively realize compute-heavy functions such as IDSes. Another limitation of switch-based middlebox implementations is that they cannot support many concurrent flows (*e.g.*, IP 5-tuple-based), due to its limited on-chip SRAM space, whereas CPU server-based implementations can leverage its large DRAM space to handle per-flow states.

Network monitoring. As programmable switches are located at a vantage point where they have high visibility of the network, various network monitoring techniques have been implemented on programmable switches [13, 78, 146, 183, 199]. To deal with limited on-chip SRAM space while supporting high accuracy, many of them leverage memory-efficient probabilistic data structures, including sketches that approximate certain characteristics of packet streams, such as heavy hitters, cardinality, and entropy. Despite the use of memory-efficient data structures, they do not scale in terms of number of concurrent monitoring dimensions (*e.g.*, IP 5-tuple, Source IP, and Destination IP) due to limited on-chip resources, including pipeline stages and ALUs. Another type of monitoring technique is in-band telemetry (INT) [13], which makes each packet carry in-network telemetry information such as switch IDs, switch queue lengths, and

processing latency on custom INT header fields. Such information could be collected by using programmable switches, and end hosts will consume it for telemetry purposes.

Improving network performance. Congestion control is one of the challenging tasks in data center networks and the Internet. While congestion control algorithms running at end hosts typically infer the congestion status (*e.g.*, via packet losses or delays) and use them to adjust sending rate, recent work has shown that with the help from programmable switches (*e.g.*, providing switch queue information via INT [13]), they can outperform existing approaches [108, 140]. Active queue management (AQM) is another type of applications enabled by programmable switches. AQM is a set of algorithms designed to shorten the queuing delay by preventing switch buffers from being full. In fixed-function switches, it is difficult or infeasible to modify or change AQM algorithms. However, recent work has demonstrated that various AQM algorithms can be implemented on programmable switches by leveraging their stateful operations as well as traffic management features and showed they could improve the end-to-end application performance [61, 201].

Security defenses. With the ability to handle Tbps of traffic at line rate, programmable switches become an appealing platform to implement volumetric network-layer DDoS attacks [148, 197, 203]. By detecting and mitigating attacks in the network, switch-based defenses can prevent end hosts from being affected by attacks cost-efficiently, compared to alternatives, such as using proprietary appliances or CPU server-based implementations. However, similar to network middlebox functions, switch-based defenses cannot prevent every type of attack or implement every type of defense due to the limited capabilities of the switches. For example, the switches cannot efficiently serve the function if a defense function requires reassembling TCP flows to detect or mitigate an attack (*e.g.*, to inspect application-layer header fields).

Accelerating data-intensive systems. Data-intensive systems such as distributed databases, key-value stores, data analytics, and machine learning systems involve a massive amount of data movements or require coordination between nodes, and CPUs often become the performance bottleneck. Recent work has shown that programmable switches can help accelerate such distributed systems by offloading components of the systems to the switches [115, 137, 138, 139, 174, 189, 200, 206]. Proposed systems reduce the amount of data that need to be processed by CPU servers by effectively performing necessary computations (*e.g.*, data aggregation, sequencing, and serving cached data) on the fly at the switches.

In the rest of this thesis, we will focus on *stateful* in-network applications developed for programmable switches, such as middlebox functions and security defenses, where correctly maintaining state is critical for their performance and correctness.

Chapter 3

Table Extension Architecture for Memory-Intensive In-Network Applications

In this chapter, we describe how to enable *elastic memory* for in-network applications using network functions (NFs) as a use case example. NFs are an essential component in today's online service infrastructure. They are deployed on the critical path of the infrastructure (*e.g.*, at the front-end) where a large volume of traffic with many concurrent flows needs to be handled. It requires NFs to be scaled for overall network operations.

NFs have been traditionally deployed either using standalone hardware appliances or a cluster of commodity servers (also known as network function virtualization (NFV)) [90, 163]. More recently, another approach has been gaining attention in the community: NFs implemented on programmable switch ASICs (*e.g.*, [13, 24, 153]).

However, we find that none of these approaches can handle NFs when there is a combination of a large number of concurrent flows (*e.g.*, $O(10M)$) and a very high traffic rate (*e.g.*, > 1 Tbps). A programmable switch ASIC cannot serve a large number of concurrent flows that requires a large flow table due to its small on-chip SRAM space although it has enough capacity to process a very high traffic rate. Similarly, it requires several tens of hardware appliances or hundreds of servers to handle the high-traffic rate, which significantly increases operational cost.

We observe that the limited on-chip SRAM space is a key bottleneck for programmable switch ASICs. If we could enable the switch ASICs to store lookup tables on cheaper DRAM in a scalable way, it could be a new enabler to serve a broader set of operating regimes which are defined by workloads and operating conditions (*i.e.*, traffic rate and the number of concurrent flows that NFs have to process), cost-effectively. In this chapter, we envision a new system architecture called TEA (Table Extension Architecture) that enables the switch ASICs on the top of racks in an NFV cluster to leverage DRAM on commodity servers.

While using server DRAM is an appealing low-cost and scalable solution, accessing server DRAM is inherently slower than accessing on-chip SRAM. As we discuss in [Section 3.2.2](#), without

careful design, this can significantly degrade processing performance and availability of NFs. Indeed there are several technical challenges in realizing this vision in practice:

- First, for external DRAM access, while RDMA (Remote Direct Memory Access) looks like a promising solution, it is unclear how to do RDMA from the switch ASIC without modifying it. Our insight is that by leveraging the programmability of ASIC, we can implement a subset of the RDMA protocol that suffices for our rack-scale deployment model in NFV clusters.
- Second, since each external DRAM access incurs high latency (a few μs), TEA must complete table lookups in a single-round trip to DRAM and must continue processing other packets. At first glance, it would seem that conventional cuckoo hashing [161] would suffice. However, cuckoo hashing is not suitable for external DRAM because it can require multiple memory accesses at times. Fortunately, we find that bounded linear probing [205], a design originally created for improving cache hit rates, can be a basis for enabling table lookups guaranteed to complete in a single round trip. In addition, we adapt this data structure to provide temporary storage to support our deferred packet processing needs.
- Third, to support NFs that require several hundred million lookups per second, we need mechanisms to leverage the available DRAM and DRAM-access bandwidth across multiple servers. While traditional distributed hashing schemes (*e.g.*, consistent hashing [119]) help scale out the lookup throughput by distributing table entries and balancing lookup request load across servers, we observe that they consume too many ASIC resources. We show that simpler, resource-efficient hashing schemes, combined with a small on-chip SRAM cache, can address both the load balancing and scaling requirements.
- Lastly, for high availability, one may detect servers' availability changes (due to server failures or congested link) in the control plane, but it could take several milliseconds to make the data plane react to it, degrading overall performance. We demonstrate that it is possible to repurpose existing ASIC's features to support rapid failure detection and fail-over in the data plane.

TEA provides a *virtual table abstraction* for lookup tables stored across the combination of on-chip SRAM and external DRAM, creating the illusion of large, high-performance tables to NFs. Our focus is on NFs such as L4 load balancers, firewalls, NATs, VXLAN or VPN gateways that are *compute-light* and *state-heavy*. Developers can write such NFs using a library of TEA APIs implemented in P4 [31] which is a programming language for programmable switches. We expose the APIs as modularized P4 codes so that developers can easily integrate TEA with their NF implementations.

We implement a prototype of TEA in P4 and four canonical NFs using the TEA API. We evaluate it with microbenchmarks as well as NF benchmarks in our testbed consisting of a Tofino-based programmable switch and 12 commodity servers. Our evaluations show that TEA allows NFs running on the switch to look up table entries with low and predictable latency (1.8–2.2 μs), and the throughput can be scaled linearly by recruiting more servers (138 million lookups per second with 8 servers in our testbed). Compared to server-based NFs with a single server, TEA-based NFs achieve up to $9.6\times$ higher throughput and $3.1\times$ lower latency without

	Hardware appliance	Commodity Server	Programmable Switch
Performance	40 Gbps	10 Gbps	3.3 Tbps
Memory	O(10GB) DRAM	O(10GB) DRAM	O(10MB) SRAM
Price	>\$40K	\$3K	\$10K
Energy consumption	480W	200W	620W

Table 3.1: Comparison of NF deployment options. We excerpt the information from product briefs [10, 28, 53] and prior work [153, 163].

consuming the CPUs and many ASIC resources. We also show that TEA can react to server availability changes within a few microseconds.

3.1 Motivation

NFs are deployed in many network settings, including inside the cloud and at the edge. They perform a wide range of tasks, ranging from packet filtering and load balancing to encryption and deep packet inspection. In this chapter, we focus on *compute-light* and *state-heavy* NFs, such as L4 load balancers, firewalls, NATs, VXLAN or VPN gateways. Even though NFs in this category are not compute intensive, they still need to support a large volume of traffic and concurrent flows on the critical path (*e.g.*, at the front-end of the cloud). Thus, their performance and scalability are the key for overall network operations.

There are three typical options to realize such NFs today: (1) using standalone hardware middlebox appliances, (2) implementing them on a cluster of commodity servers (*i.e.*, NFV cluster) [65, 90, 163], and (3) implementing them on emerging programmable switches [21, 53]. We note that while there are other options such as implementing NFs on FPGA boards attached to servers (*e.g.*, [94]), we consider the above three options that have been widely studied and deployed today.

Network operators may choose different options by considering the performance, memory size, cost, and energy efficiency of each option based on their workloads and operating conditions (*i.e.*, traffic rate and the number of concurrent flow that NF instances have to process). To understand which option is better in which scenario, we analyze a canonical NF, load balancers, in four operational regimes.¹ Table 3.1 compares these options in terms of performance, memory size, price, and energy consumption, and we use these numbers in our analysis below.

- **Regime 1: Low traffic rate (<100 Gbps) / Small number of concurrent flows (*e.g.*, 100K flows and ≈ 1 MB per-flow state).** This regime can be served by using any of three options. While supporting 100 Gbps traffic would require 3 hardware appliances (\sim \$120K), or 10 servers (\sim \$30K), a single programmable switch can support it with on-chip SRAM which is large enough to serve the small flow state. Thus, using a programmable switch would be the most cost and energy-efficient solution for this regime.

¹While our analysis focuses on a specific case of load balancers, these observations also apply to other NFs such as firewalls, gateway functions, NATs, and ACLs.

- **Regime 2: Low traffic rate (<100 Gbps) / Large number of concurrent flows (e.g., 10M flows and \approx 100 MB per-flow state).** A programmable switch cannot handle this workload since it does not have enough SRAM to store the flow state. As mentioned above, supporting 100 Gbps traffic would require 3 hardware appliances or 10 servers. In both these options, the systems can easily store the relevant flow state.
- **Regime 3: High traffic rate (>1 Tbps) / Small number of concurrent flows (e.g., 100K flow and \approx 1 MB per-flow state)** In this regime, using a programmable switch would be the most cost and energy-efficient solution because the per-flow state can fit in its SRAM space and it can easily serve the traffic. Hardware appliances and commodity servers would require many nodes to support this traffic rate making them very expensive ($25 \times$ \$40K appliances vs. $100 \times$ \$3K servers vs. $1 \times$ \$10K switch).
- **Regime 4: High traffic rate (>1 Tbps) / Large number of concurrent flows (e.g., 10M flows and \approx 100 MB per-flow state)** Many servers or appliances are required as the traffic rate increases (e.g., 10 Tbps requires 1000 high-end servers, which costs \$3M). Although programmable switches can handle the traffic rate [53], their limited memory makes it infeasible to support the needed flow state. One could add more on-chip SRAM (\$2-5K per GB) with chip modification or more switches to address the memory limitation, but costs would rise significantly.

In summary, our analysis suggests that: (1) servers and appliances can handle the low-bandwidth regime effectively, (2) programmable switches are great when flow-state fits in the limited SRAM space, and (3) nothing handles the most demanding workloads well. Ideally, if we could build an architecture that enables switches to utilize more memory with cheaper DRAM (like servers) in a scalable way, it would make programmable switches more broadly applicable and serve the extreme regime cost-efficiently.

3.2 Overview

3.2.1 Design Space

Building on the above analysis, we explore if and how we can potentially leverage external DRAM that already exists in the network. Now, there are two places where we can naturally find available DRAM *near* the switch ASIC:

- **Switch’s control plane.** The control plane has a few GB of DRAM to manage the control plane data. An ASIC could access the DRAM via the PCIe channel between the ASIC and the control plane CPU. Note that the PCIe channel has a limited bandwidth which is lower than the ASIC’s per-port bandwidth. While this low and fixed bandwidth is enough to process occasional control plane traffic, it cannot support higher traffic rates (which can cause high memory access rate) without significant hardware modifications. Also, although in theory, it is possible to add additional DRAM to the control plane, in practice, the size is fixed at design time. (e.g., 8 GB in the switch in our testbed [32]).

- **On-board off-chip DRAM.** Some switch ASIC vendors have added custom off-chip DRAM on the switch board [15]. This DRAM is used for custom tasks such as buffering packets or storing specific lookup tables. Similar to the control plane case, the memory access bandwidth and size is fixed at design time, which makes it very hard to scale without chip modification. Note that while a future switch ASIC architecture might provide on-board off-chip DRAM with larger size and higher bandwidth, it requires new interfaces and mechanisms to access DRAM from a programmable pipeline. We discuss this further in [Section 3.6](#).

We observe that two options above do not scale in terms of memory access bandwidth and capacity today, which are typically fixed at hardware design time. We believe that support for scaling becomes more critical as the total amount of traffic (both in terms of traffic volume and number of concurrent flows) each switch needs to process increases [25, 81].

Our vision. In this chapter, we take an alternative approach that leverages *DRAM in commodity servers in NFV clusters* in a scalable way. A typical NFV cluster (either inside the cloud or at the edge) consisting of multiple racks of servers [65, 90, 95, 163] already has several tens of GB of DRAM on each server. If we can reserve some portion of DRAM and let the switch ASIC located at the top-of-rack (ToR) access it, the ASIC could make use a large per-flow table, which would not be possible with on-chip SRAM today.

Using a single server could still limit the access bandwidth, *i.e.*, minimum of network bandwidth between the ASIC and the server, and PCIe bandwidth in the server. However, we can leverage multiple servers to increase the aggregate bandwidth. Also, while the ASIC uses DRAM in servers, CPUs on the servers can simultaneously serve other tasks such as compute-intensive NFs, including traffic en/decryption or payload inspection, which cannot be supported by switches today.

If this can be realized, programmable switches can become an effective way to serve *high traffic rate involving a large number of concurrent flows*, and thus work for all the regimes we considered earlier. However, realizing this vision has key design and implementation challenges, as we describe next.

3.2.2 Design Challenges

To understand why it is challenging to realize this vision, let us consider a natural starting point based on prior work using traditional Remote Procedure Call (RPC) mechanisms [117, 160] ([Figure 3.1a](#)). Specifically, the switch ASIC sends and receives RPC requests and responses via the switch control plane to avoid adding complexity (*e.g.*, state management for reliable transport) to the data plane. While this is functionally correct, there are three fundamental bottlenecks:

(1) High and unpredictable latency. A table lookup can result in high latencies because of the latency between the ASIC, the control plane CPU, and the server CPU (over the network), which can take a few hundred microseconds. Moreover, the uncertainty introduced by the scheduling logic on the switch control plane and server CPU can introduce jitter and high variability [127].

(2) Limited memory access bandwidth. The lookup throughput is constrained by the minimum of the bandwidth between ASIC-to-the-control-plane-CPU and control-plane-CPU-to-server-

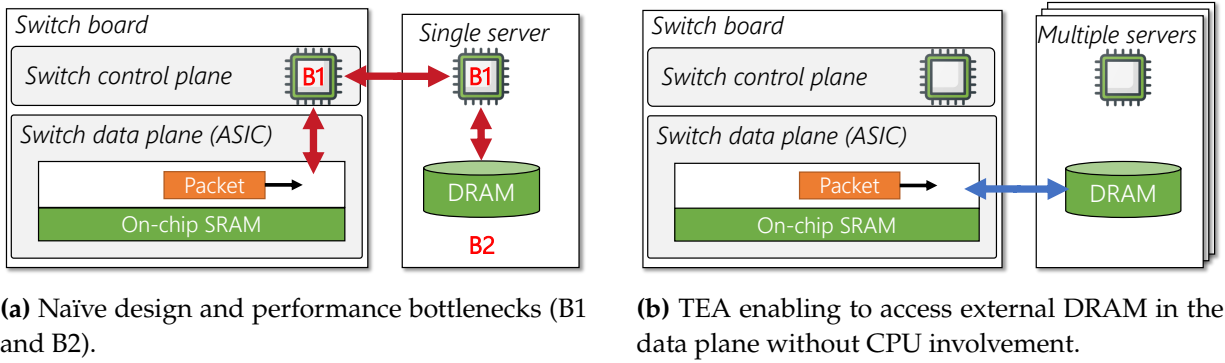


Figure 3.1: Comparison between RPC-based naïve design and TEA to access external DRAM.

CPU. Both bandwidths are typically very limited (*e.g.*, PCIe bandwidth between the ASIC and the control plane is a few tens of Gbps which is much lower than a few hundreds of Gbps of ASIC’s per-port bandwidth available today) and fixed at hardware design time.

(3) Availability. If the server fails or the network link between the control plane and the server becomes unavailable, the switch cannot lookup tables on external DRAM, degrading NF performance.

We observe that the root causes of these problems are (1) the involvement of CPUs at the control plane and the server and (2) the use of the single server (Figure 3.1a). This motivates us to ask: Is it possible to allow the switch ASIC to access external DRAM *purely in the data plane and without servers’ and the control plane’s CPU involvement* in a scalable way across multiple servers? To answer this question, we must address the following challenges:

Challenge 1. Data-Plane External DRAM Access. Switch ASICs typically do not have direct external DRAM access capability. Is it possible to enable it without hardware modifications?

Even if the ASIC can somehow directly access external DRAM, it can incur a few microseconds of latency which is an order of magnitude slower than its packet processing speed. This long latency creates the following two challenges:

Challenge 2. Single Round-Trip Table Lookups. If we use conventional hashing (*e.g.*, cuckoo hashing [161]) for storing and locating table entries in external DRAM, multiple DRAM accesses may be required to lookup an entry. Is it possible to make the ASIC do a table lookup in a single round-trip to DRAM without involving server CPUs and hardware modifications?

Challenge 3. Packet Processing. The ASIC must be able to continue processing the packet (*e.g.*, modifying header fields) after completing the lookup from external DRAM. In the meantime, it also needs to keep processing subsequent packets in the pipeline. How can we manage the packet until the lookup completes?

Challenge 4. Load-Balanced Bandwidth Use. Although using multiple servers (*i.e.*, adding network links) increases external DRAM access bandwidth, a subset of links could become overloaded due to the access locality (*i.e.*, most of memory accesses are destined to the subset of servers’ DRAM). This makes it hard to utilize available link bandwidth. How can we ensure that memory access loads are balanced across servers?

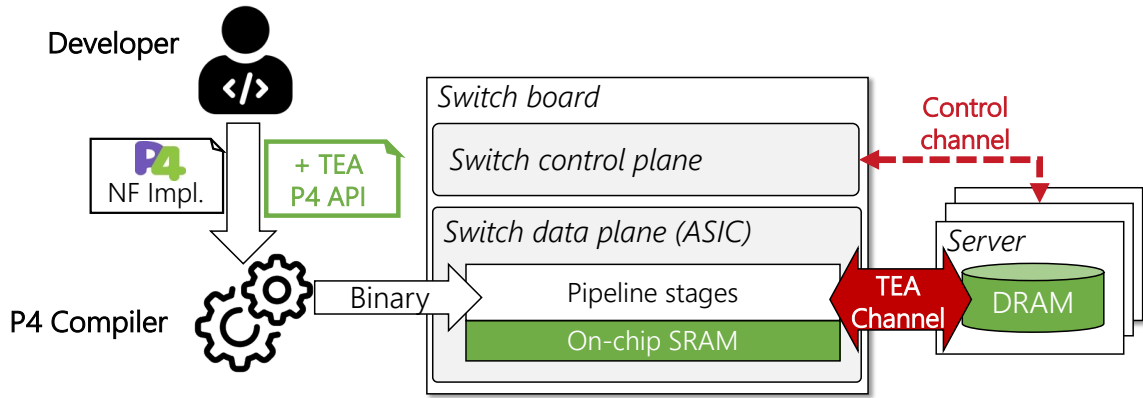


Figure 3.2: NFs implemented in P4 can be extended with TEA P4 API to look up tables across external DRAM and on-chip SRAM. The control plane is (dotted lines) involved when establishing a TEA channel.

Challenge 5. Tolerating Server Churn. Access to external DRAM becomes unavailable when a server fails or the network becomes congested (causing packet drops). How can we detect and react to these events quickly to minimize performance degradation?

3.3 TEA Design

To address the above challenges, we design TEA, a virtual table abstraction for tables stored across local SRAM and external DRAM. Using the abstraction, NFs running on a ToR programmable switch can perform key-based (*e.g.*, 5-tuple of an IP packet) table lookups, and TEA fetches the corresponding entries either from switch-local SRAM or remote DRAM. When it accesses DRAM, it delays the processing of the packet corresponding to the lookup request without blocking the rest of the packet processing pipeline. TEA’s lookup response handler resumes the delayed packet’s processing when DRAM lookup completes.

Figure 3.2 illustrates this workflow. TEA provides a set of APIs implemented in P4, a language to program NFs on programmable switches, and exposes each component as a module in P4 [31, §13]. This enables developers to easily integrate TEA with their NF implementations in P4. Once developers write their NFs using TEA components, the unmodified P4 compiler generates a binary of TEA-enabled NFs that can be loaded to the data plane and control plane APIs that can be used for configuring TEA components in the data plane.

TEA builds on the following five key ideas to address the challenges described in Section 3.2.2:

1. *Leveraging ASIC programmability* to enable simplified RDMA in the data plane (Section 3.3.1).
2. *Repurposing bounded linear probing* to guarantee hash table lookups in a single-round trip to external DRAM (Section 3.3.2).
3. *Offloading packet store to external DRAM* to enable asynchronous lookups (Section 3.3.2).
4. *Leveraging the small-cache theory* [91] to scale out the throughput (Section 3.3.3).

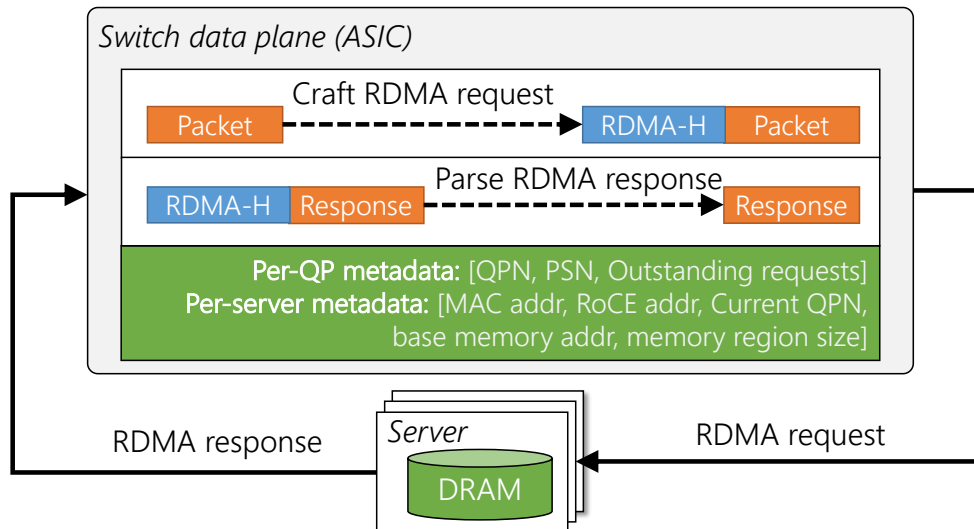


Figure 3.3: Switch ASIC generates RDMA requests by adding RoCE headers on incoming packets and parse RDMA responses without specialized capabilities for RDMA. To maintain reliable channels, the ASIC maintains per-QP and per-server metadata.

5. Repurposing ASIC's hardware capabilities to detect and react to sever availability changes in the data plane (Section 3.3.4).

3.3.1 DRAM Access in the Data Plane

To access external DRAM, we choose RDMA, which is quite common in service provider deployments [104, 155]. In comparison to RPC, RDMA is an attractive option because it is designed specifically for predictable performance memory access. It provides hardware support for a set of low-level memory operations such as read, write, and a few atomic operations (*e.g.*, fetch-and-add). Since it does not involve the server CPU for either the memory access or the reliable transport of messages, RDMA reduces both memory access latency down to $\approx 2 \mu s$, and delay jitter, and allows the use of the CPU for other compute-intensive tasks.

Challenges of using RDMA from switch ASICs. However, we still need to address two practical problems: (1) Is it feasible to generate RDMA packets purely in the switch data plane when DRAM access is needed? (2) Can we support reliable RDMA transport within the switch data plane? (*i.e.*, can switch ASICs maintain the necessary per-connection RDMA context and protocols?)

Our approach. While it may be hard to implement reliable RDMA in general on a programmable switch, we observe that we do not need fully functional RDMA for our use case. Our key insight here is that the programmable features of modern switch ASICs together with the scoped deployment model of TEA enable us to implement a small but sufficient subset of RDMA features we need.

(1) Generating RDMA packets. With respect to the first sub challenge we note that the most popular RDMA technology today is RoCE (RDMA over Converged Ethernet) protocol [113, 114],

where RDMA requests and responses are regular Ethernet packets with RoCE headers. This means that ASICs can generate valid RDMA requests by crafting RoCE packets without needing any RDMA-specific hardware components.

Figure 3.3 illustrates this high-level idea. When the data plane needs to access DRAM, it crafts an appropriate RDMA packet by adding a series of specific RoCE headers to the incoming packet. This include Ethernet headers, global route headers, base transport headers, and RDMA extended transport header with RDMA metadata such as a queue-pair number (QPN), a packet sequence number (PSN), a remote access key (Rkey), a remote memory address, and a length of data to be written or read from the DRAM.² The needed metadata is provided via the control plane in advance.

(2) Reliable RDMA. To address the second question of reliable RDMA, we leverage the assumption that in TEA, DRAM servers are directly connected to the ToR switch. This means that if we can make RDMA request and response packet not be dropped at the switch or NICs, the RDMA channel becomes reliable. Thus, we can simplify the RoCE protocol with two possible options. One is by ensuring the underlying Ethernet network is lossless via Priority Flow Control [6]. In this option, a NIC sends a PAUSE request to the switch when RDMA requests are buffered more than its threshold to prevent packet drops due to buffer overflow. When the switch receives a PAUSE request, it has to buffer packets until the NIC allows to send packets. We adopt this option in our prototype implementation in addition to our simple switch-side flow control to cope with the current NIC configuration as we describe in Section 3.4.3.³ Alternatively, we can also configure a higher QoS-level for our RDMA traffic over lossy fabric [49]. These options allow us to enable RDMA between the ASIC and DRAM servers with a minimal amount of RDMA context metadata and without complex retransmission schemes. Specifically, it only needs to maintain a QPN (4 bytes) and tracks a packet sequence number (4 bytes) and the number of outstanding requests (2 bytes) for each queue-pair, which are used when crafting RDMA requests for the QP. Maintaining such metadata in the data plane requires only up to a few KBs of SRAM in total.

3.3.2 TEA-Table: Lookup Table Structure

The design of TEA's table data structure, TEA-Table, addresses two key issues: (1) how to complete a lookup in a single round-trip to external DRAM and (2) how to defer processing of the current packet until the lookup completes and continue processing other packets without blocking. TEA-Table repurposes a data structure that was originally designed for improving cache hit rates in software switches [205] to achieve single RTT lookups and incorporates remote packet buffers within the data structure to accommodate deferred packet processing.

²QP is the connection abstraction used in RDMA communications (similar to the socket) and QPN is a unique identifier assigned for each QP. RKey is assigned to each memory protection domain where allocated memory region is registered.

³In our experiments, we observe that our switch-side flow control mechanism prevents a NIC buffer from being overflowed before the NIC generates PAUSE frames.

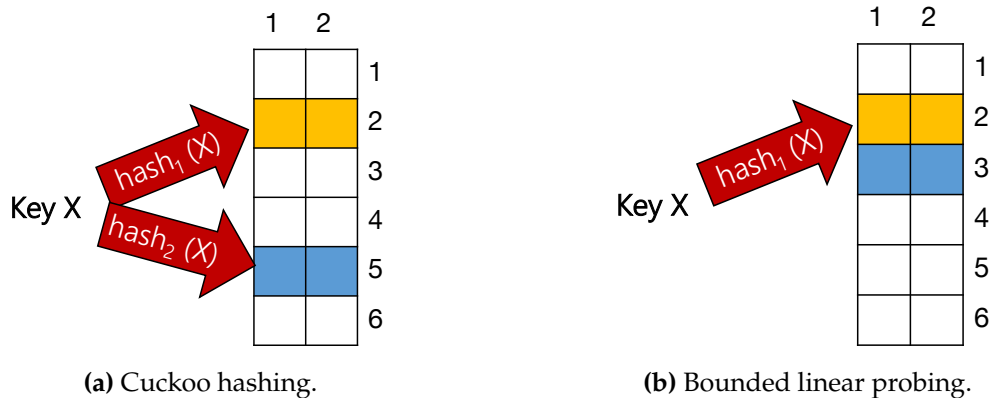


Figure 3.4: Cuckoo hashing and bounded linear probing. In this example, there are 6 buckets and 2 cells per bucket. The numbers on the top and right side indicate cell and bucket indices respectively.

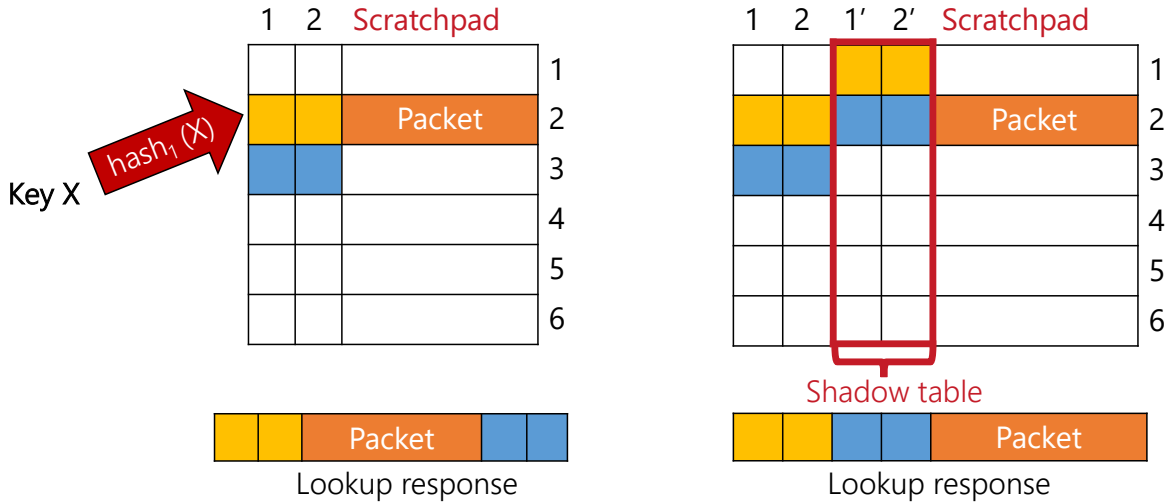
Single Round-trip Lookups:

RDMA only provides low-level memory operations such as read and write, using virtual memory addresses. However, NFs require *richer key-based lookup interface* to retrieve table entries with keys (e.g., an IP 5-tuple for an address mapping table in NAT) from DRAM. Thus, TEA must map a key to a *virtual memory address*. The challenge is that due to relatively large DRAM access latency ($\approx 2 \mu s$), we must be able to locate and fetch the entry in a single DRAM read.

Strawman solutions. At first glance, it appears we can use traditional hashing techniques. Indeed, many modern switch ASICs adopt variations of cuckoo hashing [161] for exact-match lookups in SRAM as it guarantees constant-time lookup. A caveat, however, is that each lookup requires multiple memory accesses. This means, with two-way cuckoo hashing, each lookup requires two independent memory reads. While this is feasible with fast parallel lookups on SRAM, our experience suggests that extending it to external DRAM via RDMA channel would either significantly degrade the performance of NFs or make the data plane logic complicated. To reduce multiple DRAM accesses in cuckoo hashing, we need to know precisely which of the two hash tables to access for a given key. Recent work, EMOMA [165], uses additional Bloom filters [70] in SRAM to address this issue. By checking for membership, the query can be directed to the appropriate hash table. Since there is a risk of false positives in the filter, EMOMA has a more complex item insertion that checks if inserting a new entry causes false positives. Unfortunately, this makes it impractical.⁴

Our approach. We build on a recent approach called Bounded Linear Probing (BLP) [205]. BLP was originally designed for improving cache hit rates and reducing lookup latency in software switches. Somewhat serendipitously, we find that it can also be used in our setting. [Figure 3.4](#)

⁴In our simulation, it takes several hours to insert just a few tens of million entries and implementing BFs for such a scale consumes other resources across multiple packet processing stages in the ASIC. Since such a slow insertion speed with a non-negligible amount of resource consumption makes this approach impractical, we do not consider this design.



(a) Incorrect design: switch cannot parse entries in blue cells.

(b) Corrected design with shadow table.

Figure 3.5: Design of TEA-Table with scratchpads. Scratchpads temporarily store original packets during lookups. i^{th} bucket of the shadow table has a copy of $((i + 1) \bmod n)^{\text{th}}$ bucket of the original table ($n = 6$ in this figure).

illustrates the differences between cuckoo hashing and BLP. When placing and looking up a table entry, instead of using two hash functions as in cuckoo hashing (Figure 3.4a), BLP uses one hash function and lets the second bucket be placed right next to the first bucket (Figure 3.4b).

We find that BLP's design lends itself to fetching both hash buckets in a single RDMA read. However, since BLP is designed for caching, we need to handle colliding entries differently. In BLP, when hash collisions happen, it evicts colliding entries and puts them to the main memory region (*i.e.*, DRAM). In contrast, in TEA, since the table is already located in DRAM, we put colliding entries to switch SRAM, making all entries exist in either SRAM or DRAM. Although it consumes some amount of SRAM space, we empirically prove that the collision rate is only 0.1% for the same size of the hash table as the cuckoo hash table and the same number of keys inserted. For example, when the total number of table entries is 80 million, 80K colliding entries are stored in SRAM, which takes around 4MB in the NAT mapping table with IPv6 addresses. This design is much simpler than the cuckoo hash-based approaches and requires fewer resources in the ASIC while guaranteeing at most one RDMA read per lookup.

Deferred Packet Processing:

Another key challenge is storing the packet while DRAM is accessed. This is especially critical since the $\approx 2 \mu\text{s}$ DRAM access time is very long in the context of high-speed switching where a packet is processed every nanosecond. A naïve solution would be to buffer the packet using on-chip SRAM. However, it is undesirable to use scarce SRAM for buffering a large number of packets during DRAM access.

We address this issue by storing packets to DRAM and reading back the packet along with retrieving the table entry. Specifically, we propose TEA-Table which extends our hash table structure by employing *scratchpads*. In each scratchpad, we temporarily store a packet during lookups. As shown in [Figure 3.5](#), in TEA-Table, we allocate a scratchpad for each bucket large enough to hold an MTU size packet. Note that our design requires the path MTU between the switch and the DRAM servers to be larger than the end-to-end MTU. In our prototype implementation, we set the path MTU size to 9000 bytes and the end-to-end MTU to 1500 bytes.

Hardware constraints of current RDMA NIC and switch ASIC impose another challenge. Since the NIC allows an RDMA read operation to read only a continuous memory region, with a naïve design of TEA-Table, an original packet is placed between two buckets in a lookup response, as illustrated in [Figure 3.5a](#). While we need to parse both buckets, with this format of a lookup response, the ASIC often cannot parse the second bucket (blue-colored) when the original packet (orange-colored) is large. This is because high-speed switching ASICs usually can parse only the first few hundreds of bytes in each packet.

To address this issue, we put a shadow table whose i^{th} bucket contains a copy of the $((i + 1) \bmod n)^{th}$ bucket of the original table, where n is the number of buckets in the table. As shown in [Figure 3.5b](#), the shadow table allows placing two buckets consecutively before the scratchpad in the lookup response packet. In this way, the switch can parse two buckets. Although the shadow table incurs additional DRAM consumption, given a small bucket size (<150 B) and a large available DRAM size ($>O(1$ GB)), the cost is reasonable to achieve our goal.

TEA-Table operations:

Given these building blocks, we now describe operations in TEA-Table.

- *Inserting an entry (Algorithm 1)*: Since it takes some time to complete an insertion operation, new entries are first inserted in to an *SRAM stash*, which is a small SRAM space to keep the pending entries. When there is no room in both buckets, our insertion logic running on the control plane chooses a victim cell and replaces it with the new key. In the next iteration, the logic tries to insert the key from the victim cell. If there still exists a key that fails to be inserted after *MaxTries* iterations, it remains in the SRAM Stash. Once the insertion is completed, the entry will be removed from the stash.
- *Deleting an entry*: Deletion is a simple operation which takes a key of a target entry as a parameter. To delete the entry, our deletion logic running on the control plane locates the cell of the entry using the same logic as in the insertion operation and overwrites the cell with zeros.
- *Lookup an entry (Algorithm 2)*: When an NF requests a lookup for an entry, our lookup logic first checks whether it exists in SRAM Stash or Cache (we explain the cache in [Section 3.3.3](#)), and if it does, the entry in SRAM is returned. Otherwise, after retrieving the DRAM address of the bucket, it uses RDMA to write the packet to the scratchpad of the bucket and then performs an RDMA read of the entire bucket including the packet stored in the scratchpad.
- *Lookup response handler*: Upon receiving the RDMA read request, the NIC sends an RDMA read response containing a lookup response back to the switch. To handle the lookup

Algorithm 1: Insert(key, value) for TEA-Table (Control plane).

```
1 tries=0;
2 entry=(key,value);
3 while tries < MaxTries do
    /* Temporally store the entry in SRAM during insertion */
4   insert entry to SRAM Stash;
5   i=hash(key);
6   if bucket[i] has an empty cell then
7     insert entry to the cell;
8     remove entry from SRAM;
9     copy the cell to the shadow table;
10    return Done;
11  j=(i+1) % n;
12  if bucket[j] has an empty cell then
13    insert entry to the cell;
14    remove entry from SRAM;
15    copy the cell to the shadow table;
16    return Done;
17  select a random cell c from bucket[i] ∪ bucket[j];
18  victim=c.entry;
19  insert entry to c;
20  remove entry from SRAM stash;
21  entry=victim;
22  tries++;
```

Algorithm 2: Lookup(key) for TEA-Table (Data plane).

```
1 if key exists in SRAM Stash or Cache then
2   return (SRAM[key],packet);
3 i=hash(key);
   /* Resolve memory address of the bucket */
4 addr=resolve_addr(i);
   /* Write the packet to the scratchpad */
5 RDMA_Write(addr+KV_LEN, packet, packet_length);
6 length=KV_LEN+packet_length;
   /* Read the bucket and packet */
7 (kv_cells, packet) = RDMA_Read(addr, length);
8 Lookup response handler:
9   Upon receive lookup response packet
10  return (kv_cells[key], packet);
```

response at the switch, we introduce *Lookup response handler*, which is a similar concept as the callback handler in other programming languages. Upon receiving a lookup response, the handler returns an entry and the original packet parsed from the response. TEA allows

developers to define custom actions in the handler (*e.g.*, modifying header fields with the fetched entry).

Note that as the insertion and deletion operations are relatively complex compared to the lookup operation, the control plane has to execute them. Due to this constraint, our current design does not support NFs that add and delete table entries in the data plane.

3.3.3 Multiple DRAM Servers

Recall from [Section 5.2](#), we can achieve higher lookup throughput using multiple servers. To utilize the available access bandwidth effectively, we need to answer the following questions: (1) How to partition and distribute a TEA-Table across multiple servers? (2) How to balance memory access load across the servers?

Strawman solution. To partition the table and provide load balancing, we can consider conventional distributed hashing schemes such as consistent hashing [119] and rendezvous hashing [188] as they can achieve good load balance among servers by partitioning hash tables. However, in these algorithms, each server is in charge of many non-contiguous parts (*i.e.*, buckets) of the table. In turn, this causes the switch ASIC to maintain a large number of $\langle \text{bucket range, server ID} \rangle$ mappings, consuming a non-negligible amount of TCAM space. For example, if one wants to implement consistent hashing, supporting N servers with $100N$ virtual nodes⁵ can use up to $(100N - 1)$ range-matching rules.

Our approach. Instead, we apply a simpler, resource-efficient hashing scheme to partition the table. We split the entire hash table into N sub-tables that contain buckets in a contiguous hash space and distribute them to N servers. The size of each sub-table can be different depending on the available DRAM provided by each server. This design requires only N range-matching rules in TCAM to locate a server for a key.

While this simple design reduces the TCAM usage, it may not guarantee the same load balance as the traditional distributed hashing approaches. Fortunately, we find that adding a small cache to the switch SRAM is helpful for load balancing across the servers. In particular, we leverage the theoretical results that caching at least $O(N \log N)$ popular entries where N is the number of *servers*, not the number of entries, can provide uniform load balancing across N servers regardless of traffic patterns or skewness [91]. For example, for NFs using per-flow table entries, the popularity can be defined as the number of packets in each flow. Specifically, we keep track of the popular entries within the data plane using a count-min sketch [84], for which efficient switch data plane implementations are already available [115, 147].

As an additional benefit, this cache also reduces the total DRAM access traffic in TEA. When an NF looks up the cached entries, the requests are absorbed by the switch without consuming DRAM access link bandwidth, thus reducing the number of lookup requests that need to be served by the NICs. In practice, the small cache can help achieve near switch line-rate throughput since only a few popular entries are frequently requested and consume a significant portion of throughput [67, 85, 172]. We show the effectiveness of caching for load balancing and throughput improvement in [Section 3.5.1](#).

⁵In consistent hashing, multiple virtual nodes are assigned to each physical node for better load balancing [119].

3.3.4 High Availability

As mentioned in [Section 5.2](#), TEA needs to detect and react to lookup failures to ensure high availability. We consider the following two lookup failure modes: (1) *high link utilization* due to regular network traffic (*i.e.*, other than lookup requests) could cause table lookup requests be dropped. (2) When *a server fails*, lookup requests destined to the server cannot be completed.

Strawman solution. Failures could be detected by periodically checking the port counters (to estimate link utilization) and port status (as an indicator of server failures) from the control plane. However, it could take a few tens of milliseconds from detecting an event to updating the state in the data plane. The delay can result in: (1) dropping many lookup requests due to the out-of-date state and (2) overlooking short-duration events (*e.g.*, microbursts).

Our solution. To reduce the delay, we repurpose the meter and packet generator engine of the switch ASIC to estimate port utilization and port status, respectively. Typically, the meter, which implements the RFC 2698 [110], is used for enforcing QoS policies (*e.g.*, rate limiting). When it is executed, it returns a color (red, yellow, or green) based on pre-configured rates (*i.e.*, if the utilization exceeds the rate, the meter returns red). The packet generator engine is typically configured to inject packets into a switch pipeline when a certain event happens mainly for diagnosis purposes.

To detect high port utilization, we set a threshold (link bandwidth in bps) for the per-port meter and get colors for ports where a lookup request can be routed. To detect a port down event, we configure the packet generator engine to generate a packet when ports go down. By processing the generated packet, TEA updates the port status table in the data plane. Based on these two per-port state information (utilization and status), TEA decides an egress port for a lookup request (*i.e.*, an active port that is not over-utilized). Note that since the meter is updated after a packet is completely received, it can lag behind less than a microsecond. We show that the gap is small enough to make it useful to react to high link utilization in [Section 3.5.1](#).

In our prototype, we replicate hash tables in TEA-Table to two servers and let TEA choose a server based on the availability.

3.3.5 Putting It All Together

[Figure 3.6](#) illustrates the key components of TEA on the switch data plane and servers, and how an NF uses it for packet processing. When the NF performs a lookup with a key using the TEA APIs, TEA first updates the count-min sketch of the key. Then, it checks whether an entry for the key exists in SRAM Stash or Cache (green-colored). If it exists, it directly passes the entry to the NF. Otherwise, it resolves a memory address and server ID using the memory address resolver. It then generates an RDMA write of the packet contents to the scratchpad and an RDMA read of the table row using the memory access requester (orange-colored). This design guarantees that RDMA write and read requests are always destined to the same server, and with our flow control mechanism described in [Section 3.4.3](#), both requests are not issued and a packet is dropped when the destination server is overloaded. Upon receiving an RDMA request from the switch, RDMA NICs on servers fetch entries from DRAM and send them back to the switch. Then, the lookup

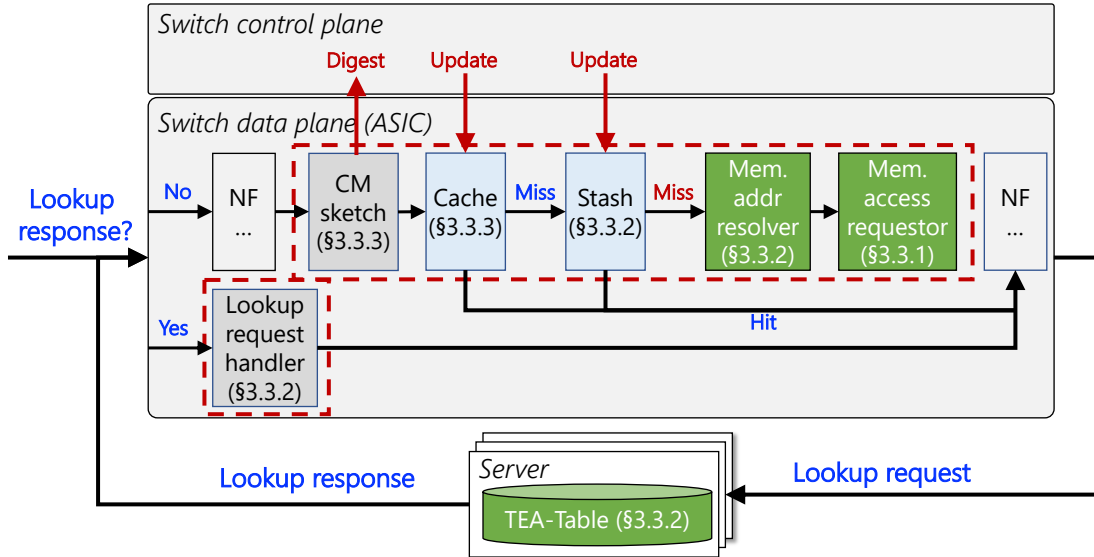


Figure 3.6: Summary of key components in TEA. The components form one logical TEA component (dotted-red box) used by an NF pipeline.

response handler extracts matched entries and the original packet contents to pass them to the NF.

Overhead of TEA. When an NF accesses external DRAM for table lookups using TEA, it incurs some amount of latency and bandwidth overheads for packet processing. For latency, as we show in Section 3.5.1, it adds up to around $2 \mu s$ per-packet latency depending on the packet size. For bandwidth, since TEA generates additional RDMA packets for external DRAM lookups, it affects both the switch pipeline and link bandwidth consumption. Within the switch pipeline, as it replicates an incoming packet to generate RDMA write and read packets, it doubles the bandwidth usage of the *egress* pipeline. It also consumes the same amount of link bandwidth between the switch and a server where a target entry is located. On the server side, while TEA does not involve CPUs, it consumes some amount of servers' memory bandwidth, which may affect performance of memory-intensive applications running on servers, especially when the memory bandwidth is fully utilized. Note that if an entry for the packet is already cached, there is no overhead.

3.4 Implementation

3.4.1 Data and Control Plane

We implement TEA's data plane in P4 [72] and compile it to Barefoot Tofino ASIC [53] with P4 Studio [45]. In the memory address resolver, we use Tofino-embedded crc64 as a hash function to locate a bucket in TEA-Table. We implement the server ID resolution using a range-matching table. In the memory access requestor, to craft lookup request packets, we make the packet replication engine in the ASIC replicate an incoming packet into two packets. The engine ensures

Network function	State	Table size (MB)
NAT	Per-flow address mapping	525
Stateful firewall	Per-flow connection state	353
Load balancer	Per-flow connection mapping	525
VPN gateway	Ext.-to-int. tunnel mapping	343

Table 3.2: The NFs we developed with TEA. Table sizes are estimated by assuming 10 million entries with IPv6 addresses.

that there is no interleaved packet between two replicas. Based on the replicas, it generates RoCE packets (*i.e.*, an RDMA write and read) by adding RoCE headers on top of the packets based on the metadata resolved by the memory address resolver.

We implement the count-min sketch [84] for collecting the statistics and determining popular entries, similar to that of prior work [115, 146]. We use 4 register arrays and 64K 16-bit slots per array to implement sketches. When the sketches detect a popular key (*i.e.*, counts of the key exceed a threshold), it reports the key to the control plane by using the digest feature in the ASIC. The digest internally maintains a Bloom filter that prevents duplicate keys from being reported. The control plane populates popular entries to the cache which is implemented as a regular exact-matching table. We use a cache of size $N=1024$ in our prototype which consumes approximately 55 KB of SRAM in NAT for IPv6 addresses.

Switch control plane and server agent. We implement the switch control plane in Python and C. It manages the ASIC via the ASIC driver using a runtime API generated by the P4 compiler. The server agent running on servers is written in C, which initializes an RDMA NICs on the servers and communicates with the switch control plane when it establishes RDMA connections.

3.4.2 Programming Network Functions with TEA

Our prototype implements TEA APIs as a library of modularized P4 codes using the concept of control block in P4 [31, §13]. Figure 5.8 shows an example program written in P4 using the TEA APIs. control block implements key modules such as the lookup response handler, memory address resolver, and memory access requestor. Extending this template, developers can integrate TEA with their NF implementations. Developers provide TEA with a definition of key (*e.g.*, 5-tuple) used of a lookup table, a structure of the table stored in DRAM (*e.g.*, using struct in C), and where to store the lookup response for further packet processing.

Based on the template, we implement NAT, stateful firewall, load balancer, and VPN gateway, described in Section 5.5 and below are the simplified P4 codes of the NFs. For example, Figure 3.8 shows how these blocks would be used to implement NAT.

To demonstrate the applicability of TEA, we implement four NFs in P4 using TEA: a NAT, a stateful firewall, a load balancer, and a VPN gateway. Table 3.2 describes the state each NF maintains using TEA and its estimated size. Brief descriptions of each are below, and simplified P4 codes are in Appendix A.

```

1  #include "tea_core.p4"
2  control Ingress (headers hdr, metadata meta) {
3      LookupHandler() lookup_handler;
4      ServerResolver() server_resolver;
5      MemResolver() mem_resolver;
6      apply {
7          lookup_handler.apply(hdr, meta);
8          if (meta.lookup_md.found == true) {
9              [Ingress NF logic]
10         } else {
11             server_resolver.apply(meta);
12             mem_resolver.apply(meta);
13         }
14     }
15 }
16 control Egress (headers hdr, metadata meta) {
17     LookupRequestor() lookup_req;
18     apply {
19         if (meta.lookup_md.found == true) {
20             [Egress NF logic]
21         } else {
22             lookup_req.apply(hdr, meta);
23         }
24     }
25 }

```

Figure 3.7: A template of P4 program using TEA abstraction. TEA exposes as a library of P4 control functions (*e.g.*, `lookup_response_handler`).

NAT. The NAT implementation uses the TEA to store NAT translation tables, to lookup a \langle private IP, Port \rangle pair for a given 5-tuple. It modifies the IP address and port header fields using lookup results.

Firewall. The firewall stores the connection state to external DRAM using TEA. For an external connection, the firewall looks up a connection state and uses it to determine how to handle packets.

Load balancer. The load balancer stores the per-flow server mapping table to external DRAM using TEA. For each incoming packet, it looks up a \langle Backend server’s IP address, Port \rangle from the table.

VPN gateway. We implement a VPN gateway (*e.g.*, [23]) based on the details described in prior work [65]. It manages the external-to-internal tunnel mapping table consisting of a \langle customer’s external tunnel ID, VM IP \rangle pair as a key and a \langle Server IP, internal tunnel ID \rangle pair as a value. For incoming packets from customers, the gateway looks up the table to retrieve corresponding server IPs and internal tunnel IDs, and translates packets.

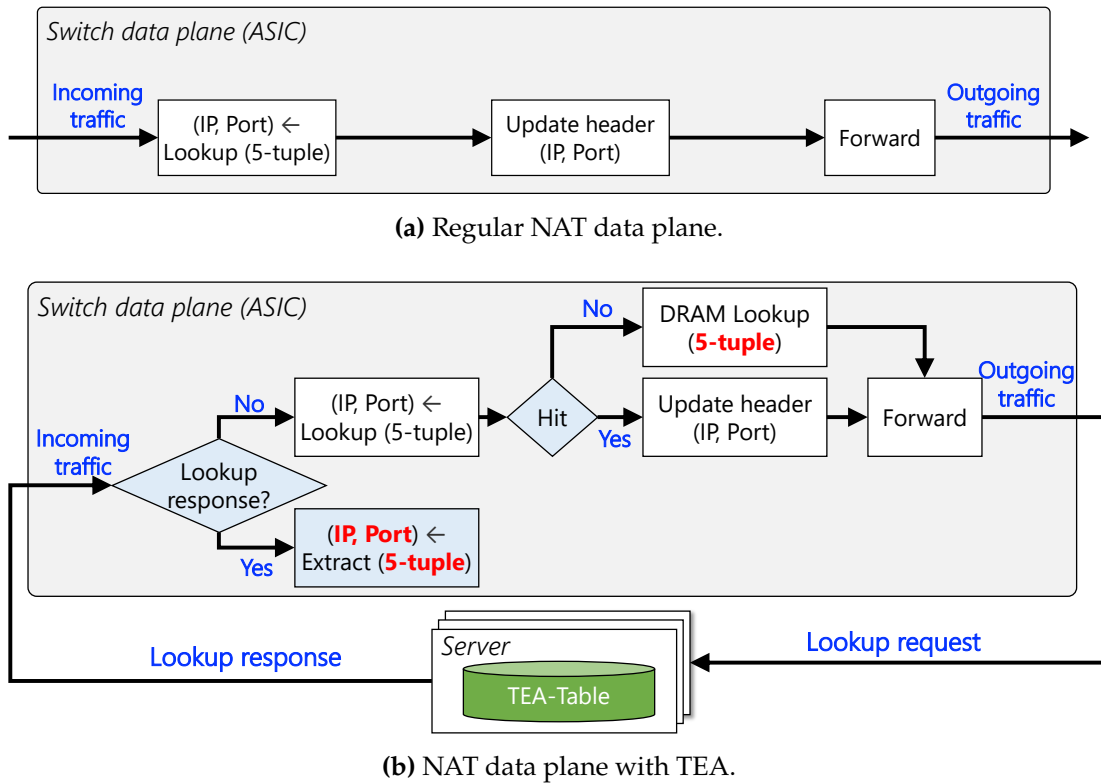


Figure 3.8: Comparison of simplified NAT data plane with and without TEA. To use TEA, in addition to the original logic (white-boxes), developers need to add TEA modules (blue-boxes) and provide basic information necessary for lookup (red-colored).

3.4.3 Limitations

The NICs in our testbed limit the maximum number of outstanding RDMA read requests to 16, and if there are more requests than the limit (*i.e.*, overloaded), they drop the requests and the QP state becomes invalid. To prevent the NICs on servers from being overloaded, we implement a simple flow control in the switch data plane, which counts and limits the number of outstanding read requests. If there is a lookup request and the number of outstanding requests has already reached to the limit, it drops the request (*i.e.*, not generating both RDMA read and write requests), causing a packet drop. This may affect the end-to-end performance. We plan to design a mechanism that routes lookup requests to an alternative DRAM server in such a case, instead of dropping packets. Also, currently, we assume that there exists at least one server that is not overloaded, and if there is no available server, TEA does not generate lookup requests and drops the packets as above.

While our NF implementations (Section 3.4.2) access one large table, some NFs may require multiple large tables. Although the current design of TEA can support multiple tables through multiple external DRAM accesses, we plan to improve its efficiency as future work.

3.5 Evaluation

We evaluate TEA on a testbed consisting of a programmable switch and commodity servers using both real data center network packet traces and synthetic packet traces. Our key findings are:

- With a single server, TEA provides a predictable lookup latency (1.8–2.2 μs) and throughput (7.3–10.9 million lookups per second) for different sizes of packets. With multiple servers, a small cache helps balance loads across servers across different skewness parameters. With the cache, adding servers scales the throughput effectively and 8 servers can perform 138 million lookups per second under a skewed workload. (Section 3.5.1).
- Compared to server-based NFs, TEA-enabled NFs are cost effective. TEA shows up to $9.6\times$ higher throughput and $3.1\times$ lower latency under the same hardware configuration. Even under an optimal setting for server-based NFs, TEA still shows $\approx 2.3\times$ higher throughput without requiring costly hardware (Section 3.5.2).
- TEA-enabled NFs can serve traffic with latency and throughput that is comparable to the switch-only implementation (*i.e.*, NFs running on a switch without accessing external DRAM) in the common case (Section 3.5.2).
- TEA provides these benefits without incurring much ASIC resource overhead. It consumes on-chip resources, including SRAM, TCAM, and hash bits, all less than 9% (Section 3.5.3).

Experimental setup. Our testbed consists of a Wedge 100BF-32X 32-ports programmable switch [32] with a Tofino ASIC and 12 servers equipped with two Intel Xeon E5-2609 CPUs (8 logical cores in total), 64 GB RAM, and a 40 Gbps Mellanox CX-3 Pro RDMA NIC. The servers run Ubuntu 18.04 with the kernel version 4.4.0. All servers are directly connected to the switch. We use 4 servers as packet generators and 8 as DRAM servers.

Traffic workloads. We use both packet traces collected from a real data center network [3] and synthetically generated ones. The packet sizes vary (64–1500 B) in the real trace. The synthetic traces are based on the observations from several data center measurement studies [67, 85, 172]. We generate packet traces with the flow size distribution in terms of the number of packets per flow that follows Zipf distribution with the skewness parameter ($\alpha=0.99, 0.95, 0.90$). We use a keyspace of 1 million randomly generated IPv4 5-tuples when creating packet traces. We generate multiple packet traces with different packet sizes and skewness parameters. We replay the traces using DPDK-pktgen [4] on packet generator nodes. In our testbed, each traffic generator node can generate 64 B packets at around 34.54 Mpps and 1500 B packets at 40 Gbps.

3.5.1 Microbenchmarks

Single-server lookup latency and throughput. First, we evaluate the performance of the DRAM access channel with a single server. For this experiment, we disable the SRAM cache. For latency, we inject 10,000 packets of different sizes (64–1500 B) to measure the lookup time. As a baseline, we setup two servers directly connected and run `ib_read_lat` in `perftest` [18] to measure RDMA read latencies for different message sizes. For throughput, we replay the trace for 30 seconds and measure the number of lookups completed during the period. Since the memory access pattern

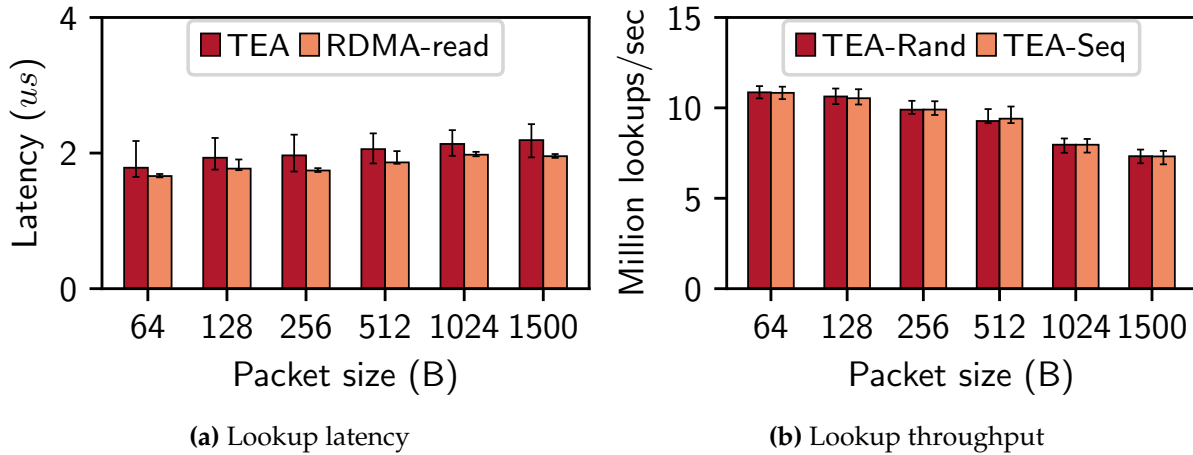


Figure 3.9: Lookup performance of TEA via an RDMA channel with a single server.

might affect the throughput, we force TEA to access buckets sequentially or randomly in this measurement.

Figure 3.9a shows the median, 10th and 90th percentile of lookup time. We see that each lookup takes 1.8–2.2 μs and the latency grows with the packet size, which is higher than raw RDMA reads (0.1–0.2 μs). This is mainly because our RDMA read request and response packets are larger than raw RDMA read packets. First, due to switch ASIC limitation, we are not able to remove an original packet from each replicated packet. This makes each RDMA read request packet have the original packet as a trailer. Second, in TEA, each RDMA read response packet consists of a bucket and the original packet, as illustrated in Figure 3.5b.

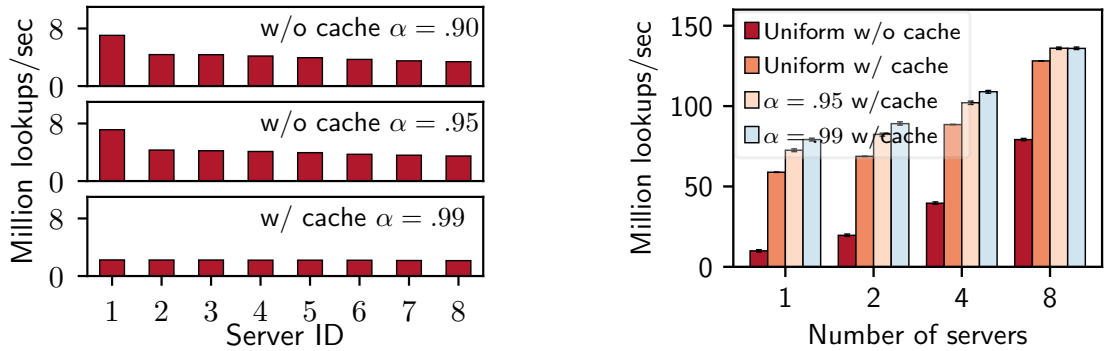
Figure 3.9b shows the lookup throughput with different packet sizes. At the maximum traffic rate we can generate in our testbed, the server NIC can handle 7.3–10.9 million lookups per second, and there is the only negligible difference (up to 0.02 million lookups per second) between sequential and random memory access patterns.

Overall, our evaluation shows TEA’s remote DRAM access channel can provide *predictable* performance which is close to the raw RDMA performance.

Throughput scaling with multiple servers. Next, we evaluate the effectiveness of using multiple servers and a small cache to scale up the lookup throughput. Here, we replay synthetic packet traces consisting of 64 B packets with the different skewness parameter (α) for the flow size distribution and measure the number of lookups served by each server with/without the cache enabled.

Figure 3.10a shows that the lookup load distribution is skewed across servers without the cache. We also observe that such a skewed access pattern limits the aggregate when the lookup request rate is high, even if there is available link bandwidth to servers. Finally, we see that with cache, even with the most skewed access pattern ($\alpha=0.99$), the load is evenly spread across servers and 49% of requests are served by the cache.

Next, we measure the aggregate lookup throughput varying the number of servers with different α values. As shown in Figure 3.10b, in all four cases, while the aggregate throughput scales linearly as we add more server, there is the difference in achievable maximum throughput



(a) Small cache balances memory access loads.

(b) Lookup throughput scales with more servers and cache.

Figure 3.10: Scalable lookup throughput of TEA with multiple servers and cache.

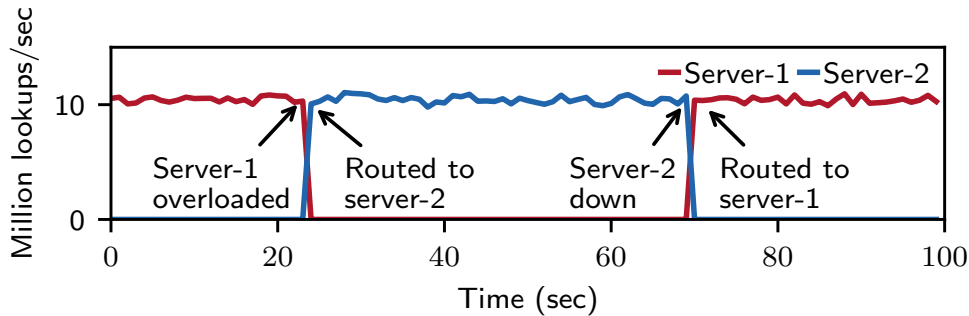


Figure 3.11: Lookup throughput changes during failover events.

depending on the skewness and the existence of the cache. We see that the more skewed the load distribution, the higher aggregate throughput TEA can support with the cache. When $\alpha=0.99$ or 0.95, TEA can process 138 million lookups per second with 8 servers and the cache. Note that this performance is limited by the maximum packet generation rate we can achieve in our testbed.

One natural question regarding the throughput would be *what is the maximum throughput an NF with TEA can achieve with N servers in a rack?* The evaluation result shows that with 8 servers TEA can support up to 138 million lookups per second. If we extrapolate this result, it means that the NF can process up to $138/8 \times N$ million packets per second, which is not high enough to support very high traffic rate with small size packets, especially when skewness is not high, and this is a limitation of our current design. For example, to support a few billion packets per second traffic rate, TEA requires more than a hundred servers, which is way more than a number of servers typically existing in a rack and a number of switch ports. Note that this analysis may not be perfectly accurate because as mentioned above, the measured maximum throughput is capped by the packet generation rate in our testbed. We plan to analyze the system throughput by injecting packets at higher rates with more servers.

Network func.	TEA w/o cache		TEA w/ cache		Server-based	
	Lat. (μs)	Tput. (Mpps)	Lat. (μs)	Tput. (Mpps)	Lat. (μs)	Tput. (Mpps)
NAT	2.34	10.64	1.93	79.37	5.62	8.49
Stateful firewall	2.35	10.58	1.91	79.23	5.59	8.37
Load balancer	2.33	10.61	1.91	79.34	5.64	8.37
VPN gateway	2.30	10.80	1.92	79.45	5.99	8.25

Table 3.3: Throughput and latency of NFs implemented using TEA with a single server and corresponding software implementations running on a single server (4 CPU cores). Note that TEA does not involve the CPU on the server.

Availability. Next, we evaluate how TEA reacts to server churn by setting up 2 servers, loading the same table entries using server-1 as a primary and server-2 as a secondary server. We replay the 64 B packet trace and measure the lookup throughput by disabling the cache. For the result in Figure 5.16, we inject the background traffic from packet generators to server-1 to emulate link utilization increase. We see that TEA starts sending lookup requests to server-2, and the throughput reaches the maximum within a second (at around 24 sec.). At this point, server-2 becomes primary. We then stop injecting the background traffic and disconnect server-2 to emulate a server failure. We can see that TEA starts routing lookup requests to server-1 as soon as it detects the event (at around 71 sec.). We observe that TEA can react to the changes in the link and server availability quickly despite a slight throughput drop at the time of failure.

3.5.2 Application Performance

Comparison with server-based NFs. We note that many factors including hardware configurations (*e.g.*, number of CPU cores) and software optimizations can affect the performance of software-based NFs. Our goal here is to show the cost benefit of TEA by comparing the performance with the same hardware configuration (*i.e.*, a server connected to a switch). For the evaluation, we implement NFs described in Table 3.2 using Click-DPDK [26] which is one of popular ways to implement high-performance NFs. We run them on the server described above.

For a fair comparison, we focus on a per-packet processing latency and throughput for 64 B packets with a single server for TEA and server-based NFs. We inject packets using 4 traffic generator nodes (max. traffic rate is ≈ 138 Mpps). Table 3.3 summarizes the results with median values for each experiment. Within each implementation option, there is no significant differences between NFs. Between TEA and server-based NFs, TEA shows up to $1.3\times$ and $9.6\times$ higher throughput, without and with the cache, respectively. For latency, TEA is up to $2.6\times$ faster without cache and $3.1\times$ faster with cache. TEA does not involve the server’s CPU at all during the experiments while server-based NFs fully utilize 4 CPU cores. Note that with more CPU cores, the server-based implementations could achieve higher throughput, ideally, close to the NIC’s raw performance (≈ 34 Mpps). Even compared to that case, TEA with cache can still

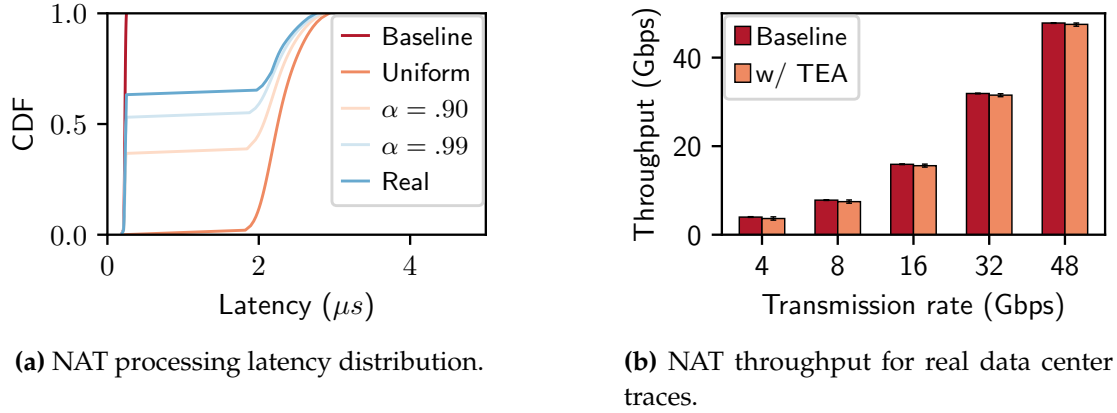


Figure 3.12: Performance of NAT using TEA.

achieve $\approx 2.3\times$ higher throughput with much lower hardware cost since it does not involve the CPU.

Comparison with switch-based NFs. To understand the overhead that TEA incurs, we compare the performance of a specific NF, NAT, running on a programmable switch, when using TEA and when using local SRAM tables (referred as baseline). The results for other NFs are similar.

To measure latency, we replay both synthetic and real data center packet traces [67] consisting of 64 B packets. Note that since the real traces consist of varying sizes of packets, we make the payload size of each packet be 64 B with the original headers (*i.e.*, the flow information is maintained). To measure the per-packet latency, we record two timestamps when packets come into the switch and leave the switch after the NAT processes the packet. Figure 3.12a shows the CDF of the latency distribution. The baseline and uniform represent the best and worst possible performance, respectively. We see that the more skewed the flow size distribution is, the lower the median latency is. Interestingly, we observe that the real traces show a skewness even higher than $\alpha=0.99$. In the traces, top 95 popular flows take more than 50% of total flows, so the cache can serve more packets, lowering the median latency. Regardless of the skewness, we see that the variance is small (no long tail), resulting in the predictable latency.

To measure throughput, we replay real data center packet traces at the rate which is higher than the original rate at which it was captured. Since the packet sizes vary, we measure the throughput in Gbps rather than Mpps. A single packet generator node can replay the trace at 14.48 Gbps, thus the maximum transmission rate we could achieve is around 57.92 Gbps with our four packet generator nodes. Figure 3.12b shows the throughput of NAT with varied transmission rates. We see that NAT with TEA can serve the traffic at the incoming rate for all cases.

3.5.3 TEA ASIC Resource Usage

We evaluate how much ASIC resource is consumed *only* by TEA based on the P4 compiler’s output. Note that as mentioned in Section 3.3.2, the number of colliding entries in TEA-Table that are stored in the SRAM is 0.1% of the total number of entries. Thus, the SRAM space usage

Resource	Additional usage
Match Crossbar	12.6%
SRAM	8.5%
TCAM	0.4%
VLIW Instruction	4.2%
Hash Bits	6.3%

Table 3.4: Additional switch ASIC resources used by TEA.

depends on the total number of inserted entries, and in this evaluation, we insert 10 million entries. Table 5.2 shows the resource consumption. We see that there are plenty of resources remaining to implement other functionality on the ASIC along with TEA. It consumes some amount of SRAM, TCAM, VLIW instruction, and hash bits, all less than 9%. Match crossbar is the most consumed resource. We observe that count-min sketch, cache, stash, and lookup response handler consume most of the match crossbar. Memory address resolver and access requestor modules consume SRAM and hash bits to store metadata for RDMA connections and resolve bucket and server IDs.

3.6 Discussion

Deployment locations. As a starting point, we focus on designing TEA for ToR switches in NFV clusters. However, TEA can be deployed in other locations. In data center racks, one can enable TEA at ToR switches with compute servers. For that, we need to make sure that there is unused DRAM space in servers and link bandwidth. Moreover, our design can be extended to non-ToR switches (*e.g.*, aggregation-layer switches) in data centers, which do not have directly connected servers under it. Since it requires multi-hop routing for lookup requests, we need to have a careful design that deals with longer and (possibly) unpredictable lookup latencies and unreliability. For example, with RoCEv2 protocol [114], which runs on top of IP/UDP and supports multi-hop routing, external DRAM access requests from upper-level switches can be routed to servers.

Match types. In this chapter, we mainly focus on exact-matching semantics. Other NFs may require other lookup types such as longest-prefix matching (LPM). Previous work emulates LPM using exact-matching [192] or converts an LPM table into a large exact-match table [126]. We can leverage such ideas to support other lookup types in TEA.

Use cases. Although the current design of TEA-Table provides a key-value based table abstraction, we can extend it to support other use cases. For example, by adopting the FIFO queue abstraction, TEA allows utilizing external DRAM as a large packet buffer which can be useful for handling packet drops due to congestion.

Other programmable switch ASICs. While we use Tofino-based programmable switches for our implementation, we believe our design can be implemented on other switch ASICs since hardware capabilities leveraged in TEA (*i.e.*, packet manipulation, meter, packet generation engine, etc.) are general features supported by most switch ASICs available today.

TEA using on-board off-chip DRAM. As mentioned in [Section 3.2.1](#), some switch ASICs support on-board off-chip DRAM for specific purposes such as packet buffers and select lookup tables [15]. As the traffic demand increases, programmable switch ASIC vendors may also consider to adopt such on-board DRAM. However, to use DRAM in a flexible manner, they need to address the same practical challenges as the ones described in this chapter, including asynchronous and low-latency DRAM access without stalling the packet processing pipeline. Thus, we believe that our techniques designed for TEA can be extended for such a future programmable switch architecture.

3.7 Related Work

Hardware-accelerated NFs. NF tasks have been accelerated using programmable switch ASICs, FPGAs, or Smart NICs to outperform CPU-only designs. Examples include offloading load balancers [153] and network monitoring [13, 105, 157] to switches and IPsec gateway, load balancer, and other NFs to FPGA-based smart NICs [94, 136]. TEA makes it possible to accelerate a wider range of NFs on programmable switches and support more operating scenarios by addressing the memory constraint issue.

Using external memory from switches. Prior work has suggested system architectures that allow switches to utilize external memory on servers [65, 122]. Such architectures run packet processing logic on both a hardware switch and a software switch on the servers and use servers' memory (*i.e.*, accessing lookup tables on servers' memory) by forwarding a subset of packets (*i.e.*, offloading traffic in certain conditions) to the software switch. This involves CPUs, increasing both average and tail packet processing latencies. In contrast, TEA purely uses DRAM on servers without involving CPUs via RDMA while addressing practical challenges in using multiple servers.

NFV state management. Previous work on state management for stateful NFs in NFV utilizes the local or remote storage to manage NF state [97, 117, 170, 196]. For example, statelessNF [117] allows NFs to leverage a centralized storage to store and load states for NFs. Their focus is better scaling and failure handling in the NFV context. In contrast, TEA leverages external DRAM to enable state-heavy NFs on programmable switches.

Other applications on programmable switches. Recent work has shown that it can be useful to offload other applications or primitives to programmable switches to enhance their performance. For example, offloading the sequencer [137], key-value cache [115, 147], and coordination service [116] improves the performance of distributed systems, in terms of throughput, scalability, and load balancing. Such systems also suffer due to switch memory constraints. TEA-like techniques could help such applications as well.

Accessing remote memory via RDMA. RDMA has been used in applications such as key-value stores [88, 118, 154], distributed shared-memory [88], transactional systems [79, 89, 127], and distributed NVM systems [150, 178]. This thesis demonstrates a novel use of RDMA, which allows a programmable switch to leverage external DRAM on such servers.

3.8 Summary

While emerging programmable switch ASIC designs make it possible for moving NFs from commodity servers to switches, the limited memory on these ASICs has been a significant impediment in their use for many NFs. To address this issue, in this chapter, we envisioned a new system architecture, called Table Extension Architecture (TEA), for top-of-rack switch ASICs in NFV clusters. TEA provides a performant virtual table abstraction for NFs on programmable switches so that they can make use of DRAM on servers connected to the switch in a cost-efficient and scalable manner. Our evaluation with microbenchmarks and NF implementations showed that TEA can provide NFs with low and predictable latency and scalable throughput for table lookups without servers' CPU involvement. Looking forward, even though our specific focus in this chapter was on NFs, we believe that TEA can be a key enabler for many innovative memory-intensive applications running on programmable switches.

Chapter 4

On-Rack Switch Resource Augmentation to Support Multiple Concurrent In-Network Applications

In the previous chapter, we show that TEA provides elastic memory via the virtual table abstraction, which is optimized for a single stateful application that uses a single large table. However, we observe that the elastic memory will not be enough to support evolving in-network application workloads. The growing popularity of in-network computing is associated with two trends: (1) the increasing number of in-network applications [109, 125] which can co-exist on a switch, and (2) the increasing demand of these applications to handle heavier workloads in terms of traffic volume and flows [25, 81]. Unfortunately, current switch resources are limited (*e.g.*, 10s MB of SRAM) and cannot keep up with the ever-increasing demands.

In this chapter, we explore an *on-rack switch resource augmentation* architecture that consists of a programmable switch and a few other data plane devices connected to the switch on the same rack. These external devices (*e.g.*, includes smart NICs [16, 41, 47, 56] and software switches running on servers [36, 190]), offer more resources to offload packet processing, albeit with some performance penalty. Perhaps more significantly, they offer a path to affordably and incrementally scale the effective capacity of a programmable network.

To effectively realize this vision of on-rack switch resource augmentation, we need the equivalent of an *operating system* to manage resources spread across multiple on-rack devices. To borrow from Anderson *et al.* [63], we can draw a first-principles analogy to the three roles that any OS serves: (1) “glue” to provide a set of common services that facilitate the sharing of resources among applications; (2) an “illusionist” to provide an abstraction of physical hardware to simplify application design; and (3) “referee” for managing resources shared between multiple applications. While there is some work on mapping a single switch app to heterogeneous devices or to augment memory (*e.g.*, [96, 128, 129, 185]), these fundamentally do not tackle multiple concurrent applications or provide these capabilities.

However, on-rack switch resource augmentation creates new challenges different from those in traditional OSEs with respect to these roles. In designing ExoPlane,¹ an OS for switch resource augmentation, we address these as follows:

- Runtime service (the glue): To avoid frequent inter-device communications during packet processing, we propose a *packet-pinning* operating model that guarantees that a packet is processed entirely on a single device.
- State abstraction (the illusionist): To enable correct stateful processing of packets even under dynamically changing workloads, we design a *two-phase state management* that places application states correctly on different devices as workload changes. We also design appropriate levels of consistency for different types of stateful objects that appear in applications.
- Resource allocation (the referee): To achieve performance and policy goals specified by developers and administrators, we formulate and solve an optimal resource allocation problem that accommodates heterogeneity across applications and data plane device capabilities.

ExoPlane consists of two key components called the planner and runtime environment. The ExoPlane planner takes multiple P4 [72] applications written for a switch with no or little modifications and optimally allocates resources to each app based on inputs from a network administrator and developers. It requires developers to add application-specific logic using our APIs *only if* the application contains a data-plane updatable object. Then, the ExoPlane runtime environment executes workloads across the switch and external devices by correctly managing state, balancing loads across devices, and handling device failures.

We implement the planner in C++, the data plane of the runtime environment in P4, and the control plane component of it in Python and C++. We evaluate it using various P4 programs in our testbed consisting of a Tofino-based programmable switch [32] and four servers equipped with Netronome Agilio CX smart NICs [16]. Our evaluations show that ExoPlane provides predictable latency (*e.g.*, ≈ 300 ns at the switch and 5.5 μ s at an external device in steady-state) and scalable throughput with more external devices (*e.g.*, up to 394 Gbps, the maximum rate in our testbed). In case of an external device failure, ExoPlane can recover an end-to-end TCP throughput within 200 ms using alternative device. ExoPlane achieves these with small control plane (a few tens MB) and switch ASIC resource overheads (less 4.5% of ASIC resources).

4.1 Motivation

In this section, we provide a primer on in-switch applications and motivate the need for resource augmentation.

4.1.1 Primer on Stateful In-Switch Applications

Many in-switch applications are *stateful*; *i.e.*, state on the switch determines how to process packets. A typical program (p) contains one or more *stateful objects* (o_i), each of which can

¹The name denotes an external (exo-) data plane

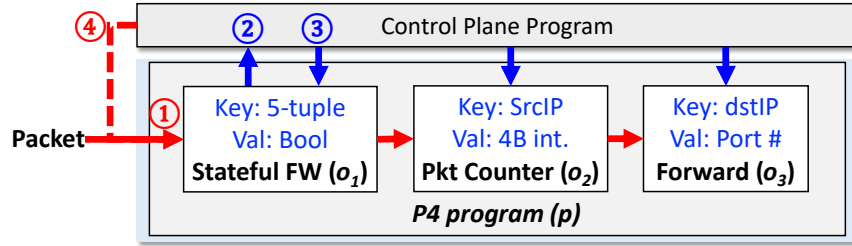


Figure 4.1: An abstract P4 application and runtime model. An application consists of multiple stateful objects (white boxes) and the control plane logic (blue arrows).

be represented as a P4 construct [72] such as a match-action table and a register.² Each object contains not only *state data* in the form of key-value pairs $((K_{o_i}, V_{o_i}))$ but also *actions*. For example, a register in P4 consists of a data array and actions that access the array. Figure 4.1 shows an example stateful P4 program (p) with three objects (o_1 – o_3). Each object requires some amount of memory (e.g., SRAM) for state data and compute resources (e.g., stateful ALUs (SALUs)) for actions. The vendor-provided compiler (e.g., Tofino compiler) allocates resources to each object using proprietary heuristics; if it cannot find a feasible allocation, the compilation will fail.

Once the program is successfully compiled and loaded to a pipeline, it can process incoming packets using its stateful objects; e.g., the firewall application in Figure 4.1 tracks active connections and drops unwanted packets from the Internet that do not belong to active connections. At runtime, the control plane logic can access the objects in the data plane (e.g., inserting a new entry to the stateful firewall (FW) object). Note that in the current switch architecture, inserting and deleting entries from a match-action table can be done only via the control plane. From the data plane, a packet only can look up an entry from the table. Registers can be read and updated by both the data and control plane. For example, in Figure 4.1, when a packet from an internal network comes in and if a state miss occurs at the stateful FW (①), it reports the packet to the control plane program (②) that inserts new entries for the packet (or flow) (③). Optionally, it sends the packet back to the data plane (④) so that it can be processed with the inserted entries.

4.1.2 Motivation

As in-network computing becomes more popular, we observe two key trends that result in increasing demand on switch resources. First, the number of applications that administrators need to run concurrently will likely increase [109, 125]. Second, the per-application workload size in terms of traffic volume and the number of flows also keeps growing [25, 81]. Unfortunately, the switch cannot keep up with this ever-increasing workloads with its limited on-chip resources (e.g., 10s MB of SRAM).

As a concrete example, suppose a cloud operator wants to deploy four applications in Table 4.1 on a network frontend switch processing traffic entering/leaving the network. Each application maintains per-flow states for each tenant to enable virtual private networks (*VPN gateway*), route

²While our focus of this chapter is on P4, other programming languages for programmable switches such as NPL [43] provide similar constructs.

Applications	States
Per-tenant VPN gateway + Packet counter	External-to-internal tunnel mapping and processed packet counter for each tenant.
Per-tenant NAT	Per-flow address mapping for each tenant. Per-flow address mapping for each tenant.
Per-tenant ACL + Filtered packet counter	Per-flow ACL and dropped packet counter for each tenant.
Sketch-based monitor	UnivMon [146] for remaining traffic classes.

Table 4.1: P4 applications deployed in a front-end switch of the data center in our motivating scenario.

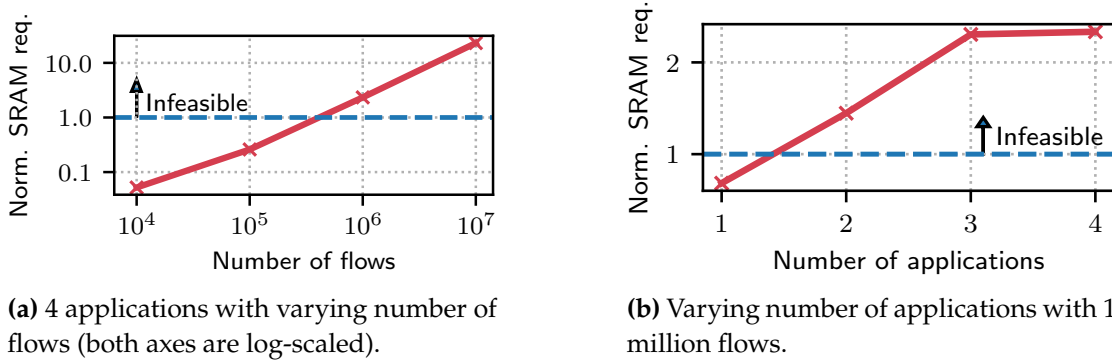


Figure 4.2: SRAM requirements (normalized to the total amount of SRAM on a switch) with varying workload size and number of applications. If the requirement exceeds 1, it is an infeasible case.

traffic from tenants’ on-premise networks to VMs running services (*NAT*), or control access to services running on tenants’ VMs (*ACL*). The *sketch-based monitor* collects statistics for the remaining traffic classes using an UnivMon sketch [146]. To see if/how these applications can coexist, we implement these applications in P4 or adopt source codes from the original authors, compose them into a single P4 program using our merger (described in Section 4.5) and compile it using the Tofino P4 compiler.

Unfortunately, we find that enabling these applications concurrently in a switch is infeasible for typical cloud workloads where an application running on a switch should be able to support at least 1M flows [85, 153, 162], as shown in Figure 4.2. We consider two scenarios: (a) running all 4 applications but varying number of concurrent flows per-application and (b) fixing number of flows to 1M but adding applications incrementally. Here, we use SRAM requirements from each application, normalized to the total amount of SRAM on a switch,³ which is the bottleneck resource in our scenario. In Figure 4.2a, we see that as the workload increases, it becomes infeasible to run all the applications. Similarly, in Figure 4.2b, we see that the switch can support

³We use normalized numbers due to NDA.

only a single application. Note that in [Figure 4.2b](#), adding the 4th application (sketch-based monitor) does not increase the SRAM usage much because its SRAM usage does not increase in proportion to the number of flows.

4.1.3 A Case for On-Rack Switch Resource Augmentation

Given these above trends, one can consider several candidate approaches; *e.g.*, optimizing applications to reduce resource footprint or adding more resources to the switch ASIC. While these are valid approaches, they have limitations; *e.g.*, applications, even if optimized, may have high resource usage, and extending switching hardware is expensive.

We explore a different practical alternative, and envision an *on-rack switch resource augmentation* architecture consisting of a programmable switch connected to a few other programmable *external* data plane devices on the same rack could be a potential solution. For example, we can assign 2U in a rack, where a programmable switch is located, to install a server equipped with four 100 Gbps Smart NICs connected to the switch. Since it only consumes a small amount of rack space and does not require any changes in other parts of the network, it provides a practical deployment model.

Note that this architecture is well aligned with technology trends. First, there are many efforts to enable P4 frontends for many data plane devices, including NPU or FPGA-based smart NICs [16, 19, 112, 193] and software switches on x86 servers [36, 44, 190]. While these devices provide lower packet processing throughput (up to a few 100s Gbps), compared to hardware switches (a few tens Tbps), they have more resources (*e.g.*, a few GB of DRAM *vs.* a few 10s MB of SRAM) to support more demanding workloads. For example, Netronome’s Agilio CX Smart NICs [16] are equipped with 2 GB of DRAM that is enough for maintaining several million flow states while being able to sustain up to 40 Gbps of packet processing rate. Second, and perhaps more importantly, we can *augment* the resources as needed by simply adding more devices as needed. Taking these above two factors into account suggests that if we could effectively realize such an architecture, it offers a cost-efficient, incrementally extensible, and potentially “future-proof” way forward to support the growing demands of multiple in-switch applications by adapting more and newer generation of external devices.

4.2 Overview

While the vision of on-rack switch resource augmentation is promising, to realize it in practice, we will need a practical OS to effectively share resources spread on multiple devices across multiple applications. Drawing an analogy from traditional computing [63], ideally this OS provides an *infinite switch resource abstraction* that gives an illusion of large resources to each application. That is, application developers and network administrators can express their programs and requirements at a higher level of abstraction without having to worry about the complexities of managing and multiplexing the resources on heterogeneous devices. While some early efforts have leveraged resources on heterogeneous data plane devices for individual in-switch applications [96, 129, 185], they do not provide the OS-like capabilities and abstractions we need for multiple concurrent applications.

4.2.1 Choice of Operating Model

A classical OS multiplexes multiple programs on the limited CPU/memory by choosing when and what processes to swap in/out. In our context, the workload is a set of incoming packets mapped to various in-switch applications, and we need to analogously choose a runtime model to multiplex the processing of these packets across applications and different devices.

Strawman models. The design space for the operating model can be defined by two dimensions: (1) Can an application run on multiple devices? and (2) Can an individual packet be processed on multiple devices? To understand the trade-offs involved, let us consider two candidate options.

- A first option is an *app-pinning* model where an application is pinned to a single device, and a packet is entirely processed on that device. In this model, since the packet is processed entirely on a single device without requiring additional logic, there is no additional processing latency due to inter-device rerouting and resource overhead. However, since the app can only run on that particular device, its throughput and available resources are limited.
- Alternatively, we can consider a *full-disaggregation* model where an application can run on multiple devices, and a packet also can be processed on multiple devices. Since an application can be placed any device, it has more available resources. However, depending on the availability of state, a packet needs to be routed between the switch and the external device multiple times. Such frequent inter-device routing increases packet processing latency and makes it unpredictable. Also, it incurs high resource overhead due to per-object inter-device processing logic to route packets to a particular device and resume processing at that object on the device. It also consumes additional link and device bandwidth.

A case for packet-pinning. We adopt a middle ground called a *packet-pinning* model that pins a given packet to one device (*i.e.*, the switch or an external device) where it is completely processed while still providing flexibility of placing an application and its flows on any devices, based on two key insights. First, it can avoid frequent per-packet inter-device routing with much lower complexity and resource overhead. Second, we observe from packet traces captured from real networks that flow key distribution is highly skewed, and only a small fraction of popular keys serves the majority of the traffic for an application (*e.g.*, 6% of keys takes more than $\approx 80\%$ of an Internet backbone traffic [57]). Thus, in this model, if we could place popular flow state entries on the switch, it allows processing the majority of traffic (*i.e.*, packets) for the application entirely at the switch while the rest of them are processed at the external device by trading off performance for more resources.

4.2.2 ExoPlane Architecture

Our ExoPlane OS implements the packet-pinning operating model via two key components (Figure 4.3):

- **ExoPlane planner** takes inputs from developers and the network administrator and allocates resources on the switch and external devices to each application.

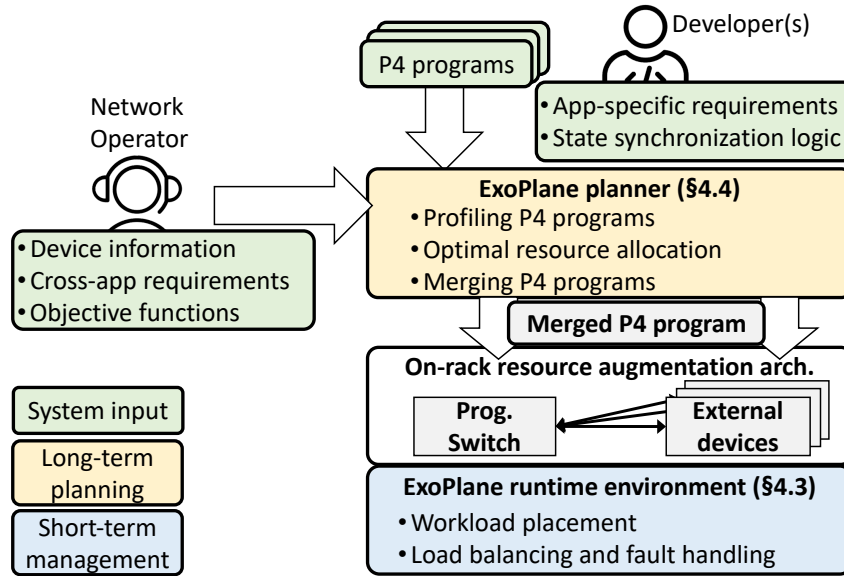


Figure 4.3: Overview: Green boxes represent the inputs for ExoPlane, and the yellow and blue box indicate key components in ExoPlane.

- **ExoPlane runtime environment** places workloads on devices, manages application states, and handles external device failures. In particular, at runtime, it tracks workload changes (*i.e.*, new flows arrive or flow popularity changes) and (re)places state object entries to the switch and external devices according to the changes.

In designing ExoPlane, we consider the following deployment capabilities: (1) A switch and external devices located on the same rack are programmable in P4-16 [31] with the same set of P4 constructs (*e.g.*, tables and registers), and we have access to vendor-provided blackbox P4 compilers; (2) External devices are equipped with large enough memory (*e.g.*, a few GB) to store entire state for multiple application. We acknowledge that not every P4-programmable device supports all the features provided by switching ASICs. According to our conversation with vendors, they plan to add such missing features, so this is not a fundamental limitation. Nonetheless, we design ExoPlane to adapt such devices as well by considering the application to device compatibility; (3) Each application handles a non-overlapping subset of traffic, which we call a traffic class so that there is no dependency between different applications (*i.e.*, a given packet is processed by only a single application);⁴ and (4) Data-plane updatable stateful objects maintain mergeable statistical data (*e.g.*, packet counter) that do not impact the control flow.

End-to-end view. As illustrated in Figure 4.3, to run applications on ExoPlane, developers provides P4 program codes and app-specific requirements (*e.g.*, affinity to the switch). Note that ExoPlane requires application modifications only if it contains data-plane updatable object whose copies can exist on multiple devices. The administrator provides information on devices (*e.g.*, resources types), cross-application workload (*e.g.*, traffic distribution), and an objective function. Then, the ExoPlane planner profiles the applications to get a compatibility to each

⁴If needed, we can apply prior offline preprocessing steps to convert overlapping subsets into an equivalent non-overlapping set [166].

device type and estimated resource footprint and performance, computes an optimal resource allocation, and generates a merged P4 program. Then, it compiles the merged program using vendor-provided P4 compilers and loads the binaries to the switch and external devices. At runtime, the ExoPlane runtime environment executes the workload (*i.e.*, packets) across the switch and external devices.

4.2.3 Design Challenges

Challenge 1. Correctness even under new flow arrivals or flow popularity changes. When the traffic workload changes, we need to (re)place object entries at the switch. We find that this can lead to incorrect packet processing due to the slow control plane operations. Also, when there are multiple copies of a data-plane updatable object on multiple devices, each of them can be updated simultaneously. While we need to synchronize them properly, we observe that it is infeasible to apply existing shared object management schemes used in server-based systems [97, 170, 196] due to hardware constraints.

Challenge 2. Handling multiple devices and device failures. While one can add more external devices to extend resources or processing capacity, we find that just adding more devices would not be effective due to possible access load imbalance across the external devices. Also, when an external device fails, we need to detect and react to the failure rapidly.

Challenge 3. Meeting objectives across applications. Given multiple applications, we have to share resources among them properly while considering per-app and cross-app objectives provided by an administrator and developers.

4.3 ExoPlane Runtime Environment

In this section, we discuss the design of the ExoPlane runtime environment. For clarity of exposition, we start with a few simplifying assumptions— a single instance of external device, steady state traffic with no workload changes, no data plane-updatable state, no device failure, and a single application. We revisit and relax these assumptions subsequently.

4.3.1 Packet-pinning Operating Model

Recall from [Section 4.2.1](#) that we adopt the packet-pinning model that ensures that each packet is completely processed at a single device (*i.e.*, requires at most a single round-trip between the switch and an external device). Here, we load an application binary and all state entries on an external device with a subset of entries loaded along with the application on the switch. As mentioned in [Section 4.1.3](#), an external device has a few GB of DRAM, which is enough to store all state entries (requiring up to a few hundred MB for a few million entries). In this model, if there is no entry for an incoming packet at the switch, the packet is routed to an external device where the packet can be processed as all state entries needed to process the packet will be presented.

However, naively implementing the packet-pinning model has two potential problems ([Figure 4.4](#)). First, if we do not carefully choose which entries to place on the switch, a high volume of

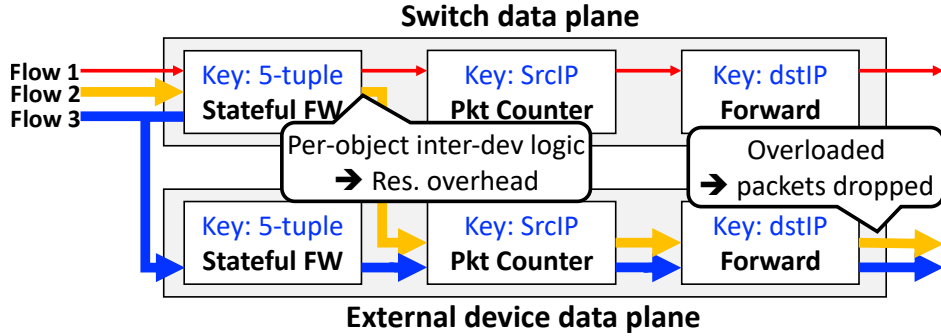


Figure 4.4: Inefficient state placement can lead an external device be overloaded, and per-object runtime components incurs high resource overhead. The arrow widths indicate relative flow volumes.

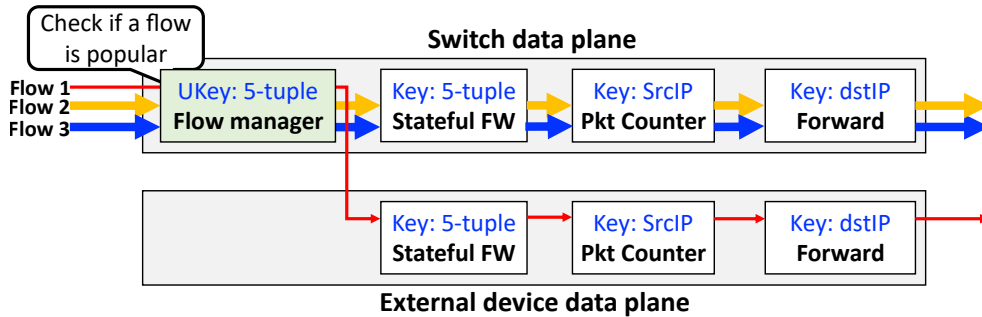


Figure 4.5: ExoPlane runtime environment processes the majority of traffic at the switch and the remaining at the external device in a run-to-completion manner. The green box is a per-app ExoPlane flow manager, and *UKey* indicates a union key of the application.

traffic will be routed to the external device (flow 2 and 3 in the figure), and it becomes overloaded, limiting the throughput. Second, since an entry miss can happen for an arbitrary object (*e.g.*, while looking up an entry of the counter object for flow 2), per-object inter-device processing logic is needed to handle such cases. Such additional logic incurs switch data plane resource overheads.

To address these problems, we propose a *union-key based* state management that enables us to process a majority of traffic for an application at the switch and the remaining at the external device (Figure 4.5). We define a union key type (*UK*) of an application as the union of key types of its constituent objects (*i.e.*, $UK = \cup_i K_{o_i}$). A flow then is a set of packets with the same union key value. For example, in the figure, an IP 5-tuple is the union key type and packets with the same IP 5-tuple forms a flow.

Having defined the union key, we can use traffic workload characteristics to enable the switch to serve the majority of traffic for the application. Specifically, we build on the observation that the distribution of flow keys (including the union key) is highly skewed in typical networking workloads. As an example, we measure the distribution of IP 5-tuple which is the union key of our example application, by analyzing packet traces collected from an Internet backbone [57]

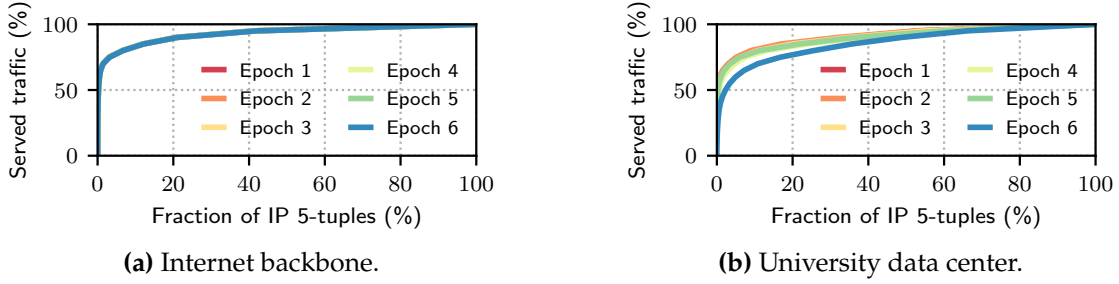


Figure 4.6: Skewness in flow key (IP 5-tuple): For both Internet backbone and data center case, a few popular keys serve the most of the traffic. This is consistent across measurement epochs.

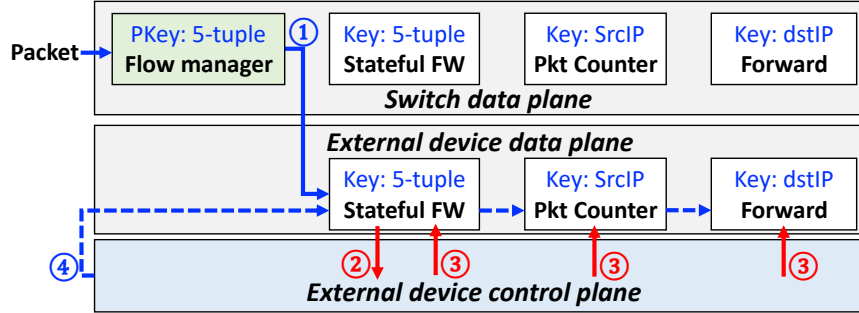


Figure 4.7: Basic workflow for inserting new flow entries in ExoPlane.

and a university data center [67]. Figure 4.6 shows the results. For both cases, we see that a small fraction of the keys contribute to the majority of the traffic; $\approx 6\%$ of keys in the backbone and $\approx 10\%$ of keys in the data center takes more than $\approx 80\%$ of traffic. The skew persists across measurement epochs (5 mins and 1 mins for the backbone and data center, respectively). We also confirm the skew exists for other coarse-grained keys such as the source IP. This suggests that we can serve the majority of the traffic at the switch by placing a few state entries with popular union keys (e.g., 516 entries for 80% in the data center trace).

Based on this, we employ a per-app *ExoPlane flow manager* (the green box in Figure 4.5 and denoted as o_{FM}) at the switch, which maintains a list of popular union keys for an application and checks whether the key of an incoming packet exists in the list when it arrives at the switch. If the key exists (i.e., the packet is from a popular flow), the packet is processed entirely at the switch. Otherwise, it is routed and processed at the external device. This way, we do not need to have per-object inter-device processing logic, minimizing the resource overhead.

In sum, our choice of packet-pinning model and its data plane design provides the following correctness property:

Invariant 1 (Packet-pinning model). For each application, if the *ExoPlane flow manager* (o_{FM}) has the packet’s union key ($UK(pkt)$), the constituent objects (o_i) must have entries ($K_{o_i}(pkt)$) for the packet.

$$\forall pkt : UK(pkt) \in o_{FM} \implies \forall i : K_{o_i}(pkt) \in o_i.$$

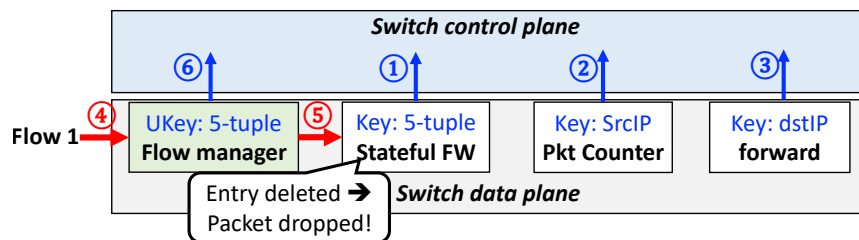


Figure 4.8: Incorrect state eviction: application’s state has been removed while there is a packet being processed.

4.3.2 Handling Workload Changes

So far we assumed steady state—(1) no new flows and (2) no changes in flow popularity. Next, we discuss how we handle new flows and then tackle popularity churn.

Handling new flows. Figure 4.7 illustrates a new flow arriving in the ExoPlane flow manager. When a packet belonging to the new flow arrives at the switch, and if a miss occurs in the ExoPlane flow manager (①), it routes the packet to the external device. Note that there are two cases for a miss happens: (1) first packet of the new flow or (2) a packet of an existing flow for which the flow state is not at the switch. Since these two cases are indistinguishable from the view of ExoPlane flow manager, it always routes packets with misses to the external device. When a packet arrives at the external device for a new flow, it must first be processed by the application’s control logic for handling new flow arrivals. In this example, the stateful FW table reports the packet to the control logic (②) that inserts entries for the flow to three objects (③). Depending on the application logic, the packet can be sent back to the data plane and processed with the new entries (④). Subsequent packets in the flow will be processed at the external device.

Promoting popular flows. In practice, the popularity of flows can change and we need to promote and demote flow states as the popularity changes. Suppose we know which flow keys become popular (*i.e.*, their entries are currently not on the switch) and unpopular (*i.e.*, their entries are currently on the switch) (We discuss how we track this in Section 4.5).

When promoting a new popular flow (*i.e.*, installing state at the switch), there are two possibilities: (1) there is spare space in the ExoPlane flow manager and application’s other objects for new entries or (2) there is no room in the objects. For (1), we can simply insert new entries to the objects. For (2), however, we need to evict some unpopular flow to make a room. As we discuss below, doing this *correctly* is challenging. Figure 4.8 illustrates why via a naïve update mechanism can violate our correctness property (1). Suppose that flow 2 becomes popular while flow 1 becomes unpopular, and there is no room for inserting new entries. Thus, the switch control plane tries to replace the entries for flow 1 with the flow 2’s. It first evicts entries for flow 1 from application objects (FW, Counter, and Forward) as well as the ExoPlane flow manager (blue arrows in Figure 4.8). However, in the current switch architecture, a set of eviction operations (blue arrows) cannot be executed atomically. Thus, there could be cases where application’s state entries have been removed already while there are packets being processed in the data plane (⑤), violating our correctness property. Even if eviction is correct, insertion can be incorrect.

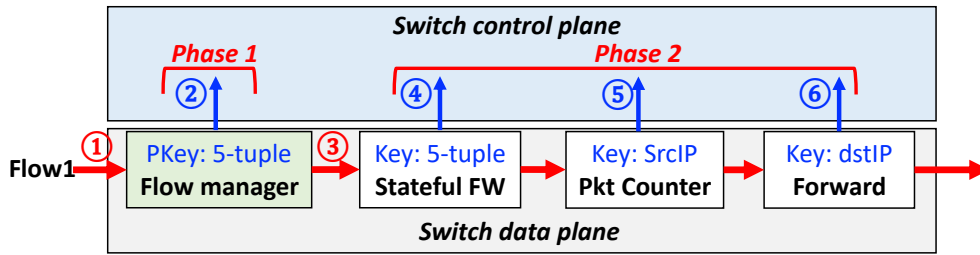


Figure 4.9: Correct two-phase state eviction.

That is, during the time the switch control plane tries to insert entries for a flow, packets for the flow arrives and are looked up the ExoPlane flow manager. If the entry exists, the packet must be processed completely at the switch. However, since entries in other objects may not be available, the packet cannot be processed and gets dropped.

Our approach: Two-phase state update. To address the issues, we propose a *two-phase state update* mechanism, inspired by classical two-phase update or commit protocols [68, 171]. As illustrated in Figure 4.9, when evicting entries for flow 1, in the first phase, the switch control plane evicts an entry from the ExoPlane flow manager. Since there can be some packets being processed in the switch data plane, it waits for a certain time period (T_{flush}) to flush out the packets. And then in the second phase, it evicts entries from the application’s objects. This mechanism ensures that all packets that arrive at the switch before the entry of the ExoPlane flow manager has been evicted are correctly processed in the switch. Note that when it evicts entries from the application’s objects, it ensures that entries for other non-victim flows will remain. The insertion works similarly. To insert entries for a flow, in the first phase, the switch control plane inserts entries to the application’s objects, and then in the second phase, it inserts an entry to the ExoPlane flow manager.

4.3.3 Synchronizing Shared Stateful Objects

The previous discussion considers scenarios with no cross-flow objects that can be updated at runtime, which meant there was no need for objects on an external device and the switch to be synchronized. In practice, applications may have some such objects; *e.g.*, per-SrcIP packet counter in our example is shared across flows. Next, we extend the basic ExoPlane protocol to handle such objects.

Consistency mode. P4 programs can have two types of stateful objects: (1) *control plane-updatable object* can be updated *only* from the control plane, such as a match-action table and (2) *data plane-updatable object* can be updated from the data plane, such as a register. Correspondingly, ExoPlane provides two levels of consistency. Control plane-updatable objects are rarely updated (*e.g.*, a stateful firewall table entry is inserted only for the first packet of each flow generated from an internal network) and an exact value is critical for correct behavior (*e.g.*, allowing packets for an established TCP connection). Thus, for this type, we provide a *strong-consistency* mode. In contrast, data plane-updatable objects can be updated more frequently (*e.g.*, per-SrcIP packet counter is updated for every packet) in the data plane and typically do not require strong

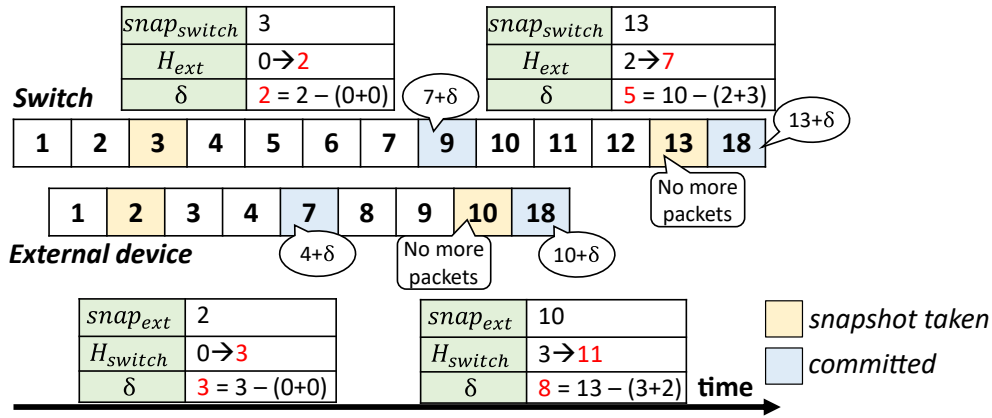


Figure 4.10: Our state synchronization protocol synchronizes two copies of an entry in the packet counter.

consistency since they maintain approximate or statistical information (e.g., packet counters and sketches). Thus, for data plane-updatable objects, we provide a *bounded-inconsistency* mode that provides consistency within a configurable time bound T_c (e.g., 1 second in our prototype).

Supporting strong consistency for control plane-updatable objects is straightforward; when the external device’s control plane receives a request for updating (or inserting) an entry to an object (② in Figure 4.7) with a key (e.g., a SrcIP), it updates (or inserts) all entries corresponding to the key existing at the external device and the switch.

Bounded-inconsistency for data-plane updatable objects is more challenging. Consider the per-SrcIP packet counter implemented using an array of registers in our example. Suppose that for a given SrcIP, there are two copies placed on the switch and the external device that can be updated simultaneously. To achieve bounded-inconsistency, the ExoPlane runtime needs to periodically *merge* values of the copies. Traditional techniques for state merging in server-based network functions (e.g., [97, 170, 196]), are impractical in our context since they rely on buffering incoming packets and pausing processing while combining copies. This is expensive and even infeasible in the switch because packet rates are much higher, and we cannot buffer packets.

Our approach for bounded-inconsistency mode. We devise a state synchronization protocol that achieves bounded-inconsistency without needing packet buffering. We do so by combining capabilities of both the switch and external device’s control and data plane. In particular, we use the control plane’s memory to track the history of periodic synchronizations while executing the merge operation in the data plane.

We explain the idea using our counter example in Figure 4.10. The control plane of each device maintains per-entry metadata including the current snapshot ($Snap$) and a history (H) of an entry value on the other side (i.e., the switch tracks the history of the external device and vice versa). For every T_c seconds, the switch control plane initiates synchronization by sending its $Snap$ and the H , and the external device’s control plane replies it with its snapshot and history; e.g., switch sends $\langle Snap=3, H=0 \rangle$ to the external device, and the external device sends $\langle Snap=2, H=0 \rangle$ back. Then, each side computes the changes that have been made at the other side (δ) after the previous synchronization by subtracting two history values from the received

snapshot value. This prevents a potential under or double-counting issue. Lastly, the control plane of both devices injects a special control packet containing δ to the data plane to combine the changes to the current state value. Note that our protocol synchronizes the copies of states correctly even when the external device fails and gets recovered. This is because the switch maintains the progress that the external device had made before the failure (H) and provides this information to the recovered device so that it can resume the synchronization from the state when it failed.

Algorithm 3: State synchronization – Switch

```

1  $S_{switch}$  : The current state of the value on the switch
2  $Ext$ : a set of external device IDs
3  $Snap_{switch}$  : The latest snapshot of the value on the switch
4  $H_{ext}[1 \dots N]$  : The latest information received from each external device
5 Upon the snapshot timer expires:
6 foreach  $ext_i \in Ext$  do
   | /* Send an initiate message to  $ext_i$  */
7   | send ( $Snap_{switch}, I_{switch}[ext_i]$ );
   | /* Receive a response from  $ext_i$  */
8   | ( $Snap_{ext_i}, I_{ext_i}$ ) = rcv ();
9 foreach  $ext_i \in Ext$  do
   | /* Adjust snapshot values and merge them */
10  |  $\delta = Snap_{ext_i} \circ^- (I_{switch}[ext_i] \circ^+ I_{ext_i})$ ;
   | /* Update the information for  $ext_i$  */
11  |  $I_{switch}[ext_i] = Snap_{ext_i} \circ^- I_{ext_i}$ 
   | /* Merge ( $\circ^+$ ) the adjusted value with the current state in the data plane */
12  $S_{switch} = S_{switch} \circ^+ \delta$ 

```

Algorithm 4: State synchronization – External device

```

1  $Snap_{ext}$  : The latest snapshot of the value on the external device
2  $S_{ext}$  : The current state of the value on the external device
3  $I_{ext}$  : The latest information received from the switch
4 Upon receiving a message from the switch ( $Snap_{switch}, I_{switch}$ ):
   | /* Send a response to the switch */
5 send ( $Snap_{ext}, I_{ext}$ );
   | /* Adjust snapshot values and merge them */
6  $\delta = Snap_{switch} \circ^- (I_{ext} \circ^+ I_{switch})$ ;
   | /* Update the history for the switch */
7  $I_{ext} = Snap_{switch} \circ^- I_{switch}$ ;
   | /* Merge the adjusted value with the current state */
8  $S_{ext} = S_{ext} \circ^+ \delta$ 

```

In [Section 4.3.3](#), we describe our state synchronization protocol to synchronize entries in a data-plane updatable object. [Algorithm 3](#) and [Algorithm 4](#) describe the detailed algorithm running on the switch and external devices, respectively.

More generally, our protocol supports objects that can be expressed in a key-value structure. We provide a developer interface to specify an object-specific *merge administrator* expressed by an addition (\circ^+) administrator that combines two values and an optional subtraction (\circ^-) administrator that subtracts one value from the other, which are used by the protocol to compute δ and commit the update. For example, a Bloom filter [70] can be expressed as (Key: an integer, Value: {0, 1}) pairs with the binary OR as \circ^+ (no subtraction administrator is needed).

4.3.4 Scaling to Multiple Devices

Thus far, we have assumed that there is a single external device. However, in practice, a single device may not provide enough processing capacity or resources, and we may need multiple external devices for this. To use multiple devices, ExoPlane shards entries in objects across the devices based on the union key. So, when an entry miss occurs at the ExoPlane flow manager, it routes a packet based on the union key to a specific external device that has state for the key. However, the skewness in the union key space (Section 4.3.1) could result in load imbalance across the devices (*i.e.*, a subset of devices can be overloaded). Fortunately, the small fraction of popular entries we already have at the switch is helpful for load balancing. Prior analysis in storage systems shows that by caching at least $O(N \log N)$ popular entries where N is the number of backend servers (in our context, external devices), guarantees uniform load balancing across the servers regardless of the skew [91]. Thus, by placing $\geq O(N \log N)$ popular union keys at the switch, we can provide the cache effect for load balancing.

4.3.5 Handling Failures

Application state loss due to failures can affect the performance or correctness of applications [130]. Specifically, we consider the failure model where an external device (or its hosting machine) fails or a network link between the switch and the device fails. To deal with state loss due to such failures, the ExoPlane runtime environment replicates each flow state to at least one additional external device when initiating entries for the flow, and when the primary device fails, it falls back to a replica. It achieves this by managing the logical to physical external device ID mapping at the switch, where the primary and replica devices share the same logical ID. However, even if there is a replica, if the runtime environment cannot detect failures and route packets to a replica quickly, the application performance can be degraded (*e.g.*, due to packet drops). To enable rapid failure detection and reaction, we repurpose the packet generation engine of the switch ASIC (which is typically used for diagnosis), similar to previous work [129]. We configure the engine to generate a packet when ports go down. By processing the generated packet, ExoPlane updates the external device ID table in the data plane. Using this, the runtime environment decides an egress port for routing the packet to a replica device.

4.4 ExoPlane Planner

Having discussed the packet-pinning for a single application, next we tackle the issue of sharing resources across multiple applications to meet the performance objectives given by developers

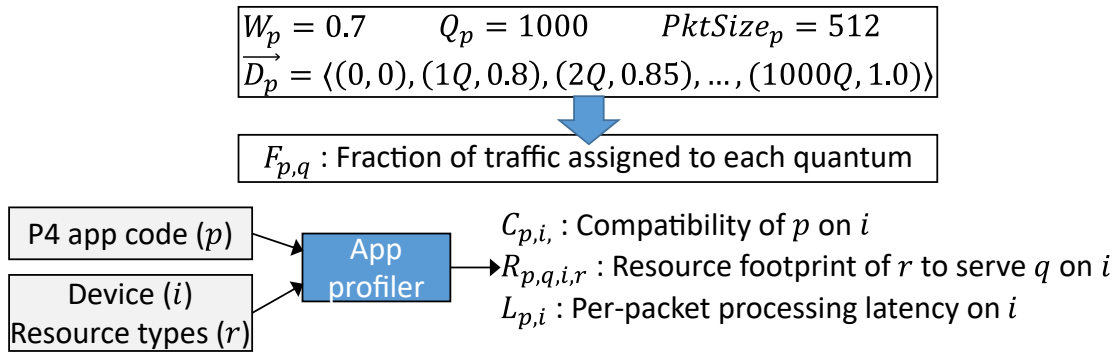


Figure 4.11: Example inputs for an application p .

and the network administrator. To this end, we design a ExoPlane planner consisting of the resource allocator and the application merger. The resource allocator takes inputs from the developers and the network administrator, finds an optimal resource allocation, and the application merger generates a merged P4 program based on the optimal allocation decision. Figure 4.11 illustrates example inputs for an application.

4.4.1 Inputs

Developers provide a set of P4 programs (p), each of which consists of a set of stateful objects. For each object, developers specify required size (*e.g.*, the number of entries in a table or register object). Optionally, they can also provide a per-app affinity to the switch as one of three values: high, medium, and low. If the affinity of an application is set to high (or low), the application will run entirely at the switch (or at external devices). The network administrator provides cross-app and per-app traffic information, which includes a fraction of traffic served by each application out of the entire traffic arriving at the switch (W_p) and the cumulative traffic distribution (D_p) over the union key space. While using a fraction of traffic served by *each key* provides the most fine-grained information, we find that it could make the search space for resource allocation too large. Instead, we use the distribution discretized into a larger quantum size denoted as Q_p . Based on D_p , we compute the estimated fraction of traffic served by each quantum q ($F_{p,q}$). The network administrator also provides resource type (r) information of device (i). For example, we consider SRAM, TCAM, hash units, and SALUs for a Tofino-based switch and compute units, SRAM, and DRAM for NPU-based NICs. The network administrator can easily extend this to other resource types.

4.4.2 Profiler

Based on the inputs, we generate per-app profiles consisting of a resource footprint, per-packet processing latency, and compatibility matrix for each device type.

The profiler estimates resource footprint of r for p serving q on i denoted as $R_{p,q,i,r}$. Since blackbox compilers determine the resource usage using proprietary heuristics, our preprocessor compiles p to obtain the usage information. For each q , it updates the size of each object specified

in an application code and compiles it using vendor-provided compilers. Then it extracts the resource usage from compiler outputs. If the compilation fails due to insufficient resources, it sets the resource usage to infinite. We use constants $Cap_{i,r}$ to represent the total amount of r available on i . The profiler also estimates a per-packet processing latency of p on i , $L_{p,i}$. Specifically, it instruments the switch to record two timestamps on a custom packet header field when a packet enters and leaves the rack. Then it injects $PktSize_p$ -sized packets to the rack and estimates the latency based on the timestamps in returned packets.

Finally, some vendor-provided P4 compilers for external devices may not support certain switch hardware features or P4 constructs (e.g., Packet recirculation and P4 registers) used by applications. Because of this, if an application uses a feature that is not supported by an external device, the device cannot run the app. To consider the compatibility of the app on devices, our profiler generates a compatibility matrix ($C_{p,i}$) that indicates whether p can be run on device i based on a set of features supported by i and a set of features used by p . The first set can be typically obtained from vendor's compiler manual. For the second set, the profiler analyzes the app code to extract used features.

4.4.3 Resource Allocation

Given these inputs, next we discuss how we formulate the problem of finding an optimal resource allocation satisfying per-app and cross-app requirements.

In our formulation, we assume that the resource usage of multiple applications can be estimated by accumulating the resource usage of each app. We use binary decision variables $d_{p,q,i}$ to indicate whether q for p is assigned to i . There are three types of constraints imposed by: (1) the assignment of q , (2) the compatibility of p on i , (3) the amount of available resources, and (4) the processing latency of p on i .

$$\forall p, q : \sum_i d_{p,q,i} = 1 \quad (4.1)$$

$$\forall p, q, i : d_{p,q,i} \leq C_{p,i} \quad (4.2)$$

$$\forall i, r : \sum_p \sum_q d_{p,q,i} \times R_{p,q,i,r} \leq Cap_{i,r} \quad (4.3)$$

$$\forall p : lat_p = \sum_q \sum_i d_{p,q,i} \times F_{p,q} \times L_{p,i} \quad (4.4)$$

First, q must be assigned to a unique i (Equation 4.1). Second, q can be assigned to i if and only if p is compatible with i (Equation 4.2). Third, the amount of r consumed by q on i must be less than or equal to the total amount of r on i (Equation 4.3). Last, the expected latency of p is the sum of per-packet processing latency of p on i weighted by $F_{p,q}$ (Equation 4.4).

The network administrator provides an objective to share resources across multiple application fairly. One possible fairness metric would be minimizing the weighted sum of the expected processing latency of each application:

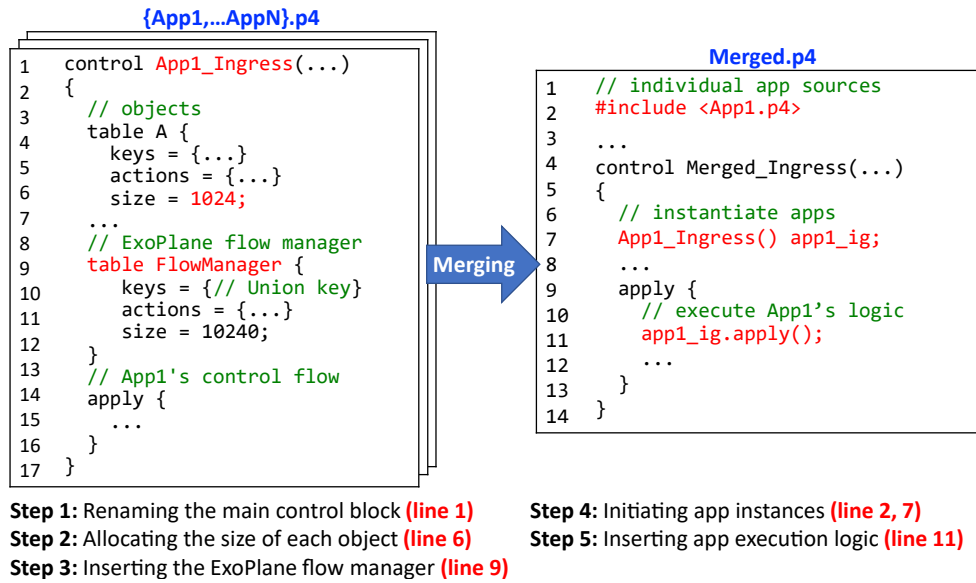


Figure 4.12: Merging multiple P4 programs into a single program.

$$\text{Minimize } \sum_p W_p \times lat_p \quad (4.5)$$

Other commonly used fairness metrics such as maximizing the minimum expected throughput can be used as well. By solving the ILP, ExoPlane resource allocator finds an optimal assignment of q to i for p , and the size of each object and ExoPlane flow manager for p accordingly, which are used as input for the application merger, as we describe next.

4.4.4 Application Merger

Given a set of P4 programs and the optimal resource allocation decision, our application merger combines the programs into a single P4 program, following our deployment model for multiple applications, described in Section 4.2.2. Our merger supports programs written in P4-16 [31]. Figure 4.12 illustrates how the merger works. First, for each application, the merger renames the main control block [31] to avoid naming conflicts between applications. Second, it specifies the size of each object (e.g., number of entries in a table) based on the decision made by our resource allocator. Third, it inserts an ExoPlane flow manager. Last, in a merged P4 code, it instantiates an instance of each application and inserts execution logic. The merged P4 code is compiled using the vendor-provided compiler and loaded to the switch and external devices. Sometimes, the compilation process could fail due to its proprietary heuristics for resource allocation. In such a case, we try more conservative allocation (i.e., with a tighter resource constraint).

In summary, ExoPlane planner allocates resources across multiple applications based on inputs from developers and network administrators and produces a merged P4 program loaded to devices. This process needs to be re-run when a set of applications or workloads changes,

which we do not expect to happen frequently (*e.g.*, for every hour). While this module is not on the critical path, we evaluate it for completeness in [Section 4.6.6](#).

4.5 Implementation

Data plane. The data plane components of the runtime environment implemented in P4-16 consists of: (1) ExoPlane flow manager implemented using a match-action table and (2) global logical to physical external device ID mapping implemented using a register array (on the switch).

Tracking flow popularity. We implement a switch resource-efficient flow popularity tracking mechanism based on the fact that in ExoPlane, packets corresponding to *potentially* popular flows are routed to an external device which has a larger amount of resources. On external devices, we use the count-min sketch [84] to track the frequently accessed flow keys. When it detects a new popular key, it reports the key to its control plane that maintains a list of reported flow keys and corresponding entries, and they are reported to the switch control plane. On the switch, we enable the *aging* supported by the switch ASIC for the ExoPlane flow manager. If a certain key of the ExoPlane flow manager has not been accessed for a timeout period (T_{idle}), a callback function registered at the switch control plane is triggered along with the information about the key. In our prototype, we set T_{idle} to 5 seconds.

Control plane. We implement the control plane components of the runtime environment in Python and C++. The main capability needed is to initialize new flow entries and promote new popular flows' entries on the switch based on the information reported by the data plane runtime components. To implement this, on the switch side, we use Barefoot Runtime APIs to access the stateful object in the switch data plane and on the smart NIC side, we use Netronome Thrift APIs [51] to interact with the NIC data plane. The switch and the external device control planes are communicated via an out-of-band TCP session over the 1 Gbps management network.

Resource allocator. We implement the resource allocator in C++ based on the Gurobi C++ API [39] to encode and solve our resource allocation ILP.

Application profiler and merger. We extend the open-source P4 compiler [46] to parse and analyze input P4 programs. Using its frontend, we extract information from each program including an entry size of each object. We implement the application merger in C++, which takes an IR generated by the compiler frontend, and produces a merged P4 program.

4.6 Evaluation

We evaluate ExoPlane on a testbed consisting of a programmable switch and servers equipped with a smart NIC using various workloads. Our key findings are:

- In steady-state, ExoPlane provides predictable per-packet latency (*e.g.*, 273–384 ns at the switch) and scalable throughput with more external devices while the app-pinning model achieves a limited throughput ([Section 4.6.1](#)).
- Even under dynamic workloads, ExoPlane can process packets with the correct state and sustain high throughput with multiple devices ([Section 4.6.2](#)).

Applications	States
Per-VM NAT	Per-flow address mapping for each VM.
Per-VM Stateful FW + Packet counter	Established TCP connection list.
Per-VM SYN proxy	Per-flow sequence number translation table.
NetCache [115]	Key-value store cache.

Table 4.2: Switch programs written in P4 used in the evaluation in addition to ones introduced in Table 4.1.

- In case of an external device failure, ExoPlane can recover an end-to-end TCP throughput within 200 *ms* (Section 4.6.4).
- ExoPlane provides the above benefits with small control plane (*e.g.*, a few tens MB) and switch ASIC resource overheads (*e.g.*, less than 4.5% of ASIC resources) (Section 4.6.5).

Testbed setup. We build an on-rack resource augmentation architecture consisting of Wedge100BF-32X Tofino-based programmable switches [32] and 4 servers equipped with Netronome Agilio CX 40 Gbps Smart NICs [16]. We use 4 additional servers to generate traffic workloads. All servers are equipped with an Intel Xeon Silver 4110 CPU and 128 GB DRAM, running Ubuntu 18.04 (kernel version 4.15.0). We repeat each experiment 100 times unless otherwise noted.

Traffic workloads. We use packet traces collected from a real data center [3], the Internet backbone [57], and synthetically generated ones. The packet sizes vary (64–1500 B) in the real trace. We generate packet traces with the flow key distribution in terms of the number of packets per flow that follows a Zipf distribution with the skewness parameters ($\alpha=0.9, 0.95, 0.99$). We use a keyspace of 1 million randomly generated IPv4 5-tuples when creating packet traces. We generate multiple packet traces with different packet sizes and skewness parameters. We replay the traces using DPDK-pktgen [4] or run iperf [14] for TCP workloads.

Deployment scenarios. We use two scenarios with multiple P4 applications: (1) at the data center front-end, 4 applications in Table 4.1 and (2) at the leaf of the network, 4 applications from Table 4.2. Given packet traces, we synthesize inputs for the resource allocator in ExoPlane planner (*e.g.*, per-app affinity and a flow key distribution). For example, we set the affinity level for the UnivMon [146] and NetCache [115] to high so that workloads for these applications are always processed at the switch.

4.6.1 Performance in Steady State

First, we evaluate the per-packet processing latency and throughput of applications running on ExoPlane in steady state (*i.e.*, no new flows, no changes in flow popularity, and no device failures). Here, we pre-populate popular flow entries at the switch and assume that the traffic is equally distributed across the applications (*i.e.*, $W_p = 0.25$ for all applications).

Per-packet processing latency. We define the packet processing latency as the time difference between when a packet first arrives at the switch from a sender and when it is sent to a receiver

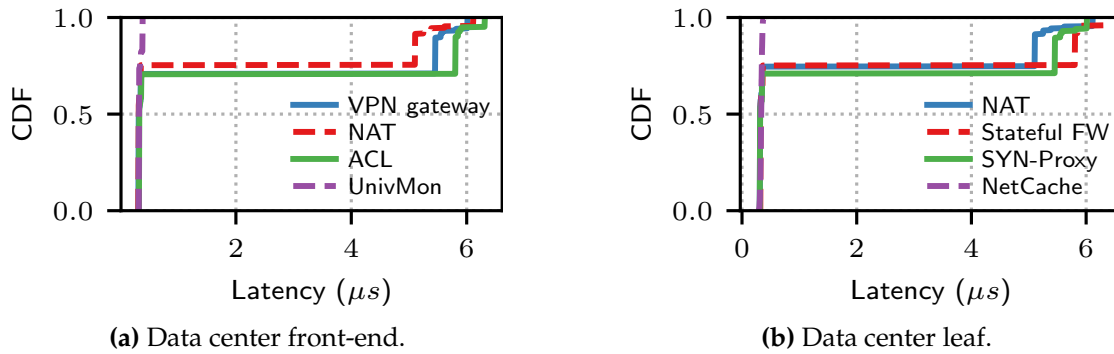


Figure 4.13: Per-packet processing latency distribution of applications concurrently running on ExoPlane in steady state.

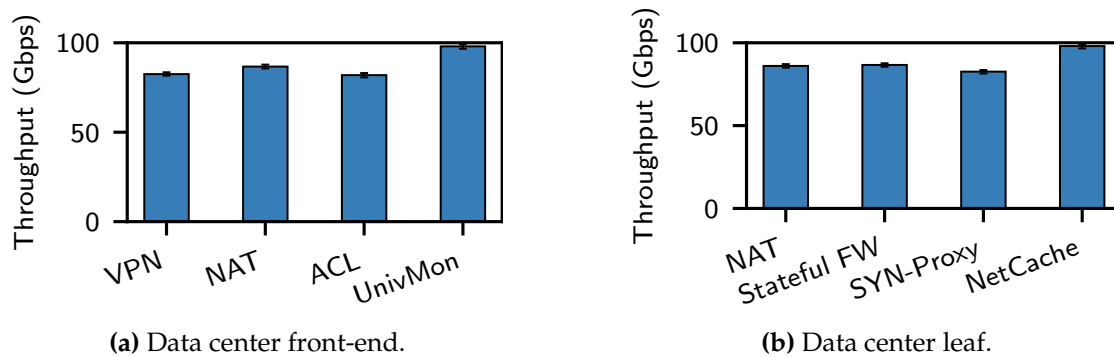


Figure 4.14: Throughput of each application running on ExoPlane in steady-state with a single external device (Applications are running concurrently).

after processing. We instrument the P4 program running on the switch to record two timestamps (48-bits each) to our custom packet header fields of each packet so that the receiver can compute the processing latency for a packet. From the sender, we replay the backbone packet traces, each of which contains more than 6 million flows.

Figure 4.13 shows the CDF of the per-packet latency distribution for each application. For the applications that are assigned to the *high* affinity (UnivMon and NetCache), every packet is processed at the switch, where each packet is processed in 273–384 ns depending on packet sizes. For other applications, the distributions vary depending on packet sizes and how much traffic is processed at the switch and the external device. The higher affinity level assigned to an application, the more traffic is processed at the switch. For example, in the front-end scenario (Figure 4.13a), at the switch, the ACL processes $\approx 70\%$ of its traffic whereas the NAT processes $\approx 75\%$ of its traffic. At the external device, packets are processed in 5.1–6.1 μs depending on an application. The key takeaway from this experiment is while there is latency gap between the switch and the external device, on each device, per-packet processing latency is predictable.

Application throughput. To measure the application throughput, we replay the synthetic packet trace that consists of 1500 B packets at line rate (98.6 Gbps in our testbed). We use four sender nodes, each of which generates traffic for each of four applications (*e.g.*, node 1 generates traffic

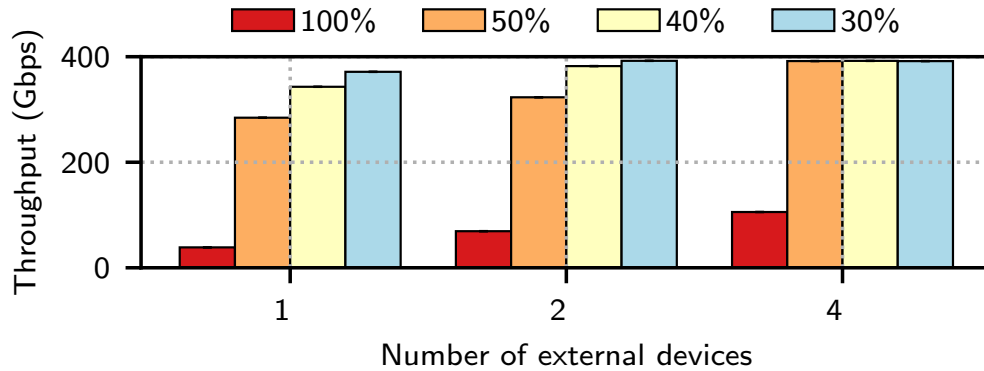


Figure 4.15: Scalable throughput with multiple devices. Each color represents a different fraction of traffic offloaded to external devices.

for application 1). Each application sends processed packets to the receiver node to measure the throughput. We start with a single external device to demonstrate the impact of the number of external devices to the throughput.

Figure 4.14 shows the throughput of each application. The applications that run entirely on the switch (UnivMon and NetCache) process traffic at line rate without dropping any packets. However, we observe that the others cannot process their traffic at line rate. This is because the aggregate amount of traffic across the applications, which needs to be processed at the external device (≈ 81 Gbps in the front-end case) exceeds the processing capacity of the single device (≈ 39 Gbps).

Scaling throughput with multiple devices. By adding more devices, ExoPlane can support higher throughput. To demonstrate this, we measure the aggregate throughput of the four applications in the front-end scenario (the maximum traffic rate is ≈ 394 Gbps) while varying the fraction of traffic offloaded to external device(s) and the number of external devices. In this experiment, we control the fraction of traffic offloaded to external devices by manually assigning the affinity of each application. UnivMon is still pinned to the switch. Figure 4.15 shows the results. In the case of 30, 40, 50% of the traffic being offloaded to external devices, we see the throughput effectively increases with more devices. In contrast, when 100% of traffic is offloaded, adding more external devices is not effective even though there is remaining capacity in the devices due to load imbalance. This results show the load balancing effect of serving popular flows at the switch, described in Section 4.3.4.

Comparison with the app-pinning model. We evaluate the benefit of ExoPlane over the app-pinning model (described in Section 4.2.1) while running 4 applications from Table 4.1. In this model, we place an application along with its entire state at the switch if there is a room. Otherwise, we place it to one of the external devices, which has the largest remaining capacity. Table 4.3 compares the aggregate throughput when running an ensemble of applications. While ExoPlane provides the maximum throughput for each ensemble, the app-pinning model achieves up to 69.3% lower throughput. This is because while ExoPlane allows an application to effectively utilize available resources across different devices, the app-pinning model fixes an application to a device.

Ensemble of applications	App-pinning	ExoPlane
VPN	98.6 Gbps	98.6 Gbps
VPN+NAT	137.1 Gbps	197.2 Gbps
VPN+NAT+ACL	174.6 Gbps	295.6 Gbps
VPN+NAT+ACL+UnivMon	271.3 Gbps	394.1 Gbps

Table 4.3: Comparison of aggregate throughput of four applications running on the app-pinning model and ExoPlane with four external devices.

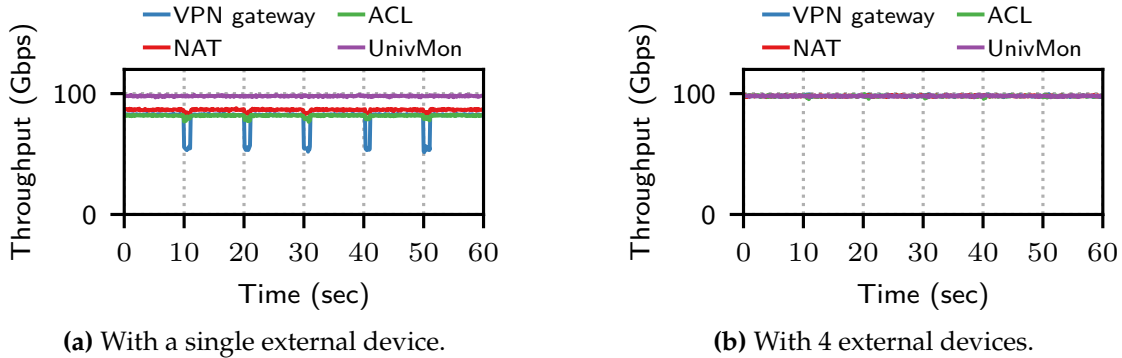


Figure 4.16: Throughput changes due to workload changes.

4.6.2 Performance under Dynamic Workload

Per-packet processing latency. As mentioned in Section 4.3.2, workload changes happen due to new flows arriving or changes in flow popularity. Handling new flows in ExoPlane can increase packet processing latency because the first packet of a new flow can be processed only after initiating necessary state for both the ExoPlane flow manager and application’s objects. In contrast, packets in the flow that becomes popular can be processed either at the switch or an external device with the same latency shown in Section 4.6.1. Thus, for each application, we measure the processing latency of the first packet of each flow. We observe that the median latency for the first packet of a new flow is 32 *ms*, which is an order of magnitude higher than that of an external device in steady state. There are two factors that contribute to this latency. First, the Neutronome Thrift API takes a few tens of *ms* to insert new entries to objects, which is not an ExoPlane-specific overhead. Second, since ExoPlane replicates entries for new flows to one another external device, it incurs additional latency when handling new flows.

Application throughput. The changes in flow popularity can impact the application throughput. To measure the throughput changes, we use the same setup as the previous measurement in steady state, but for every 10 second, we alter the most popular top 10 flows for the VPN gateway of the front-end scenario. Figure 4.16 shows the throughput changes over time. Again, we first use a single external device. As shown in Figure 4.16a, when the popularity changes, there is a sharp drop in the throughput of the VPN gateway. Also, the throughput of other applications slightly decreases as well. This is because until the state entries for the new set of popular flows

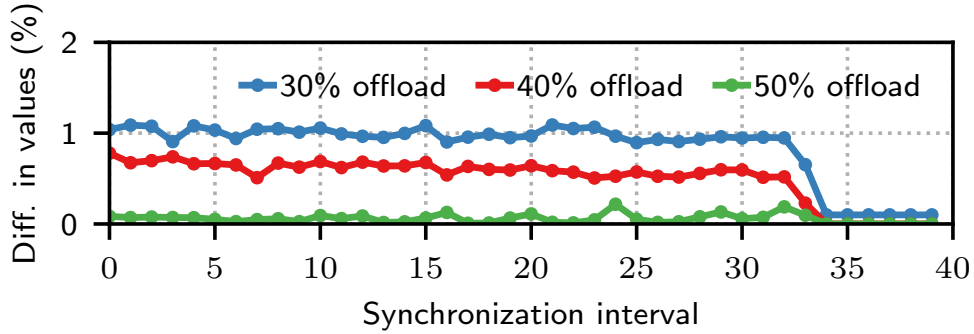


Figure 4.17: Difference in shared object values on the switch and external devices; there are no more packets after the 32nd epoch.

are installed at the switch (*i.e.*, a transient period), a high volume of traffic for the flows are routed to the external device, exceeding its processing capacity. On the other hand, as shown in Figure 4.16b, with 4 external devices, there is no such performance drop because there is enough processing capacity at the external devices to handle the traffic during the transient period.

4.6.3 Shared Stateful Object Synchronization

Next, we evaluate the effectiveness of our state synchronization protocol (Section 4.3.3) using the per-SrcIP packet counter in the stateful FW. Here, the metric of interest is the difference between the shared counter entries maintained by each device at each synchronization interval ($T_c=1$ sec.). We measure this by recording the values at each device right after executing the merge operation in the data plane while injecting 1500 B packets for 60 seconds at 98.6 Gbps. We vary the fraction of traffic offloaded to external devices. In our setting, there are 1000 entries shared between the switch and at least one of the external devices, and we get the median of the differences. Figure 4.17 shows the result with three different fractions of offloaded traffic. When the switch and external devices process the same amount of traffic (*i.e.*, 50% offload), there is almost no difference, whereas when there is a gap between the amounts of traffic (*i.e.*, 30% or 40% offload), there are some differences. This is because incoming packets keep updating the counter at each devices during the synchronization, affecting the measured values. However, we see that the variance of the difference is small across the synchronization intervals regardless of the gap, showing that our mechanism synchronizes the values. We also confirm that after the packet transmission is done, copies at each device are synchronized with the same value as the total number of packets.

4.6.4 Failover

In this experiment, we demonstrate how fast the end-to-end TCP throughput can be recovered by ExoPlane in the presence of external device failure. In Figure 5.16, we use a NAT as an example and run iperf to measure TCP throughput changes. There are 4 TCP connections,

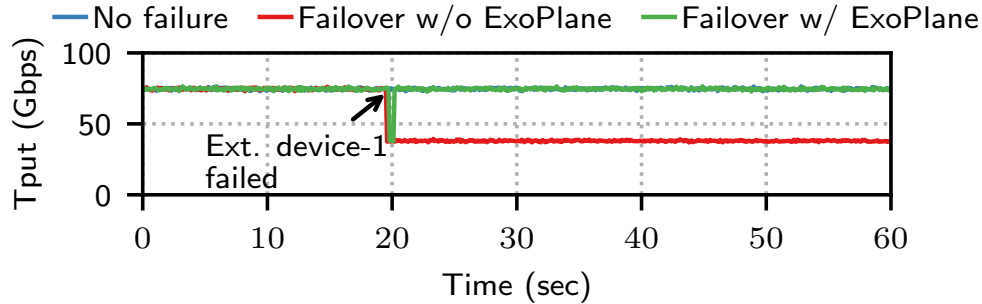


Figure 4.18: TCP throughput changes during failover and recovery.

and we configure two of them to be processed at the switch while the remaining is processed at an external device. There are two external devices enabled, and we compare changes in TCP throughput when (1) there is no failure and (2) one of the external devices fails with and without ExoPlane. We emulate the failure by disabling a port connected to the external device. At around 20 sec., when the external device-1 goes down, our failover mechanism generates a control packet that modifies the logical to physical device ID mapping in the switch data plane without involving the control plane. Then, subsequent packets are routed to the replica device. We see that the TCP throughput is recovered to its original rate within a 200 ms ⁵ whereas without ExoPlane, it cannot be recovered.

4.6.5 Runtime Resource Overheads

Control plane resource overhead. The control plane component of ExoPlane runtime environment maintains metadata for application’s states, including a mapping between union keys and devices and a history of each shared object entry on other devices. (*e.g.*, the switch control plane maintains the history of entries on external devices). Each of them consumes the control plane memory. In our scenarios, the union keys to device mapping consumes 12.5 MB per application and the history metadata consumes 1.5 MB per shared object (*e.g.*, srcIP counter in the firewall). Our state synchronization protocol consumes management network bandwidth as it periodically exchanges the information between devices, which contains a snapshot and a history of each entry. In our setting, the bandwidth consumption is 24.4 Mbps per shared object, which increases in proportion to the number of devices, the sync. interval, and the number of entries.

Switch ASIC resource usage. The data plane component of ExoPlane runtime environment consumes some switch ASIC resources. Since we implement it using an exact-match table with the aging feature and a register array, it consumes SRAM, SALUs, hash bits, MAP RAM, and match crossbar,⁶ whose usage increases proportionally to the number of popular flows maintained (except for SALUs). In our setting where 10240 popular flow entries are managed, it consumes 4.4% of the SRAM, 2.1% of SALUs, 3.5% of the hash bits, 3.8% of the MAP RAM, and 3.6% of the match crossbar, leaving ample resources to application logics.

⁵The finest sampling granularity supported by iperf is 100 ms .

⁶MAP RAMs are used for the aging feature and match crossbars are used for implementing the ‘matching’ part of match-action tables.

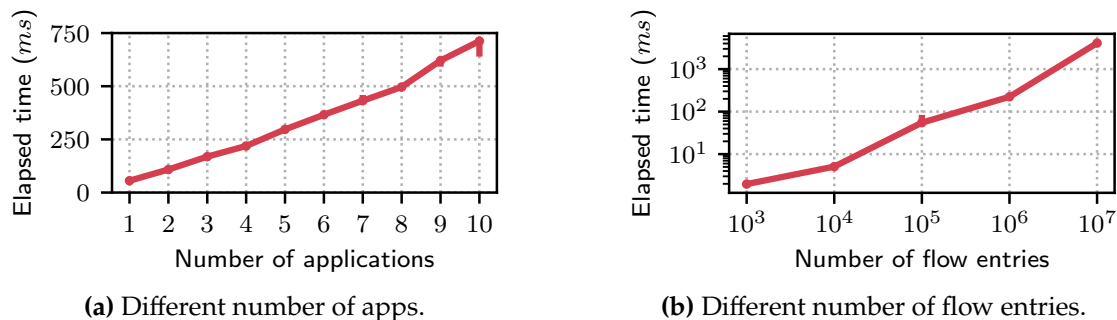


Figure 4.19: Elapsed time for the resource allocator.

4.6.6 Performance of the ExoPlane Planner

We evaluate the performance of the ExoPlane planner. In this experiment, we measure the elapsed time for finding optimal resource allocations and generating a merged P4 program on a server in our testbed. For the two sets of applications and the hardware configuration used in our evaluations, our resource allocation takes 54.5 *ms* and merging the program takes 642 *ms*, which is reasonable since the orchestrator needs to run this process on the hours or days timescale. To further understand the impact of different parameter values including the number of applications and traffic workload sizes, we synthesize inputs for the resource allocator and measure the elapsed time. First, we fix the number of external devices to 16 (to support a large number of applications) and the number of union key-based flow entries to 1 million for each application. Then, we vary the number of flow entries while fixing the number of applications to 4 and the number of devices same as the above. As illustrated in Figure 4.19a, the resource allocation time grows linearly up to 712 *ms* as the number of application increases. Also, as shown in Figure 4.19b, as the number of flow entries increases, the elapsed time also increases up to 4.1 second when each application needs to handle 10 million flow entries. This experiment illustrates the ExoPlane orchestrator takes longer time as we add more applications and increases the workload size, which can be up to a few seconds, it is still with in the reasonable timescale under our deployment model.

4.7 Related Work

Switch resource augmentation. Our work, TEA (Chapter 3) provides a virtual table abstraction that allows a single switch app to access remote DRAM for lookup tables. It is optimized for a single app with a single table. Flightplan [185] takes a single app written with custom annotations and disaggregates it to multiple devices. Developers need to be aware of external devices and manually partition the app so that each device runs only a particular portion of the app. Lyra [96] proposes a custom language for writing a single app that is disaggregated across multiple heterogeneous switches. In contrast, ExoPlane provides an OS for switch resource augmentation to support multiple concurrent P4 programs.

Language and framework for data plane composition. Some prior works attempt to support multiple data plane programs or modules in a single device [107, 184, 202, 204]. For example, virtualization approaches such as Hyper4 [107] and HyperV [202] allow composing multiple P4 programs with a constrained programming model. P4Visor [204] provides a resource-efficient merger that merges different versions of a program. However, they fail to work when the amount of resource required by the composed program exceeds the available resources in the switch and do not consider multiple heterogeneous devices.

Server-based network function state management. Previous work on state management for server-based NFs in NFV utilizes the local or remote storage to manage NF state [97, 117, 120, 170, 196]. While they work well for server-based NFs, they do not directly applicable in our setting due to workload characteristics of switch applications and hardware constraints, as pointed in [Section 4.3.3](#).

4.8 Summary

To fully unleash the potential of in-network computing, we need to look beyond the single application, fixed workload models in consideration today to support richer concurrent application workloads. Unfortunately, limited on-chip resources has been a roadblock to support multiple applications concurrently on a switch. In this chapter, we envisioned a practical approach of on-rack switch resource augmentation as an affordable and incrementally expandable solution to this dilemma. To effectively realize this architecture, we argued the need for systematic OS abstractions and addressed key challenges in realizing these abstractions. Our evaluation with various P4 applications showed that ExoPlane can provide applications with low and predictable latency, scalable throughput, and fast failover while achieving these with small resource overheads and no or little modifications on applications. Thus, ExoPlane can be a basis for a future-proof way of enabling in-network computing workloads for future apps, workloads, and emerging hardware device capabilities.

Chapter 5

Supporting Fault-Tolerance for Stateful In-Network Applications

Although with TEA and ExoPlane, we can support multiple concurrent stateful in-network applications, a key missing piece remains: fault tolerance. Classic network designs followed the end-to-end principle [173], keeping a critical state only on the end hosts. It enabled a fate-sharing approach to reliability [82]; when switches are stateless, recovering from their failure simply entails finding a new communication path. Stateful in-switch applications [12] challenge this paradigm; *e.g.*, the failure of a switch running a load balancer may cause the loss of its forwarding state, breaking thousands of active connections. While data center networks are engineered with redundant network paths [102, 172, 180] to provide fault tolerance at the routing layer, there are no capabilities for recovering in-switch states after a failure.

Thus, we need to reconsider fault tolerance for in-switch processing – something previously done in ad hoc, application-specific ways. Our goal in this chapter is to ensure that, after a failure and reroute, the same application state becomes available at the replacement switch, without degrading performance and while remaining transparent to end hosts.

Making switch state fault tolerant is uniquely challenging because of the scale and resource constraints involved. Techniques like checkpointing and active replication, which have been applied to software middleboxes [169, 179], are designed for server-based systems. These techniques rely on obtaining a consistent snapshot of state and buffering output until state updates are durably recorded to other servers. However, a switch’s high packet processing speed (a few billion packets/second [50, 53, 54]) and its limited compute and storage capabilities make it infeasible to translate these techniques to the switch context.

In this chapter, we introduce RedPlane,¹ a fault-tolerant state store for in-switch applications. RedPlane provides APIs for developers to (re)write their stateful P4 programs and make them fault-tolerant. This allows an application to retain consistent access to its state, even if the switch it runs on fails or traffic is rerouted to an alternative switch. RedPlane achieves this through a data plane centric replication mechanism that continuously replicates state updates to an external

¹The name denotes a *replicated data plane*.

state store implemented using DRAM on commodity servers. Note that running entirely in the data plane channel is key to keeping up with the switch’s full processing speed.

Realizing this high-level idea in practice entails several challenges. First, traditional notions of strict correctness with linearizability and exactly-once semantics for operations require reliable communication and output buffering. However, this is infeasible on the switch data plane due to its limited capabilities. Second, at the traffic volumes the switch data plane needs to process, naïvely requiring per-packet coordination with the server-based state store imposes severe performance overheads. Last, routing decisions when a switch fails could be unpredictable. Thus, we must be able to transparently migrate the relevant state between two switches regardless of the routing decisions.

We address these challenges with the following key ideas:

- Based on the requirements of in-switch applications, we define two practical correctness models. First, based on our observation that network applications are already resilient to packet loss, we define a strict consistency mode by explicitly adopting the standard definition of linearizability [111], which permits operations that do not complete while still providing strong consistency. Second, for write-centric applications (*e.g.*, monitoring using sketches [84]) that can tolerate approximate results, we propose a relaxed consistency mode that allows some state to be lost after a failure, but bounds the inconsistency with lower overheads.
- Instead of buffering packets using limited switch resources, we use the network itself and state store’s memory as temporary storage by piggybacking packet contents on coordination messages.
- To enable reliable state replication, we build a lightweight sequencing and retransmission protocol that ensures state updates are processed in the correct order, without requiring complex protocols (*e.g.*, TCP) in the switch data plane.
- To avoid overheads due to frequent coordination with the state store, we propose a lease-based state ownership protocol [100, 142, 158] to provide correctness without coordinating on every state access and migrate ownership between different switches as needed.

We design the RedPlane protocol that realizes our consistency modes, prove its correctness, and confirm this using a TLA+ model checker [52]. We implement a prototype of RedPlane in P4 [31] and C++ and Python, and show that different types of applications can be fault tolerant using it. We evaluate it with various applications in our testbed consisting of two Tofino-based programmable switches, four regular switches, and 10 servers. Our evaluation results show that under failure-free operation, RedPlane has negligible per-packet latency overhead for read-centric applications like NAT, and less than 8 μ s overhead even for the worst case. When a switch fails, RedPlane can recover end-to-end TCP throughput within a second by accessing the correct state.

5.1 Motivation

In-network processing has flourished in recent years, as a natural convergence of the demand for sophisticated network functionality from data center operators and the commercial availability of programmable switch platforms [21, 50, 53]. Programmable switches are used for classic middlebox functionality [123, 153], monitoring [13, 105], DDoS defense systems [197, 203] and accelerating other networked systems [115, 137, 138, 139, 147, 174, 189, 206].

These applications are *stateful*; *i.e.*, state on the switch determines how to process packets. In this chapter, we focus primarily on *hard state* applications, where a loss of state disrupts network or application functionality.² An example is an in-switch NAT, where the key state is an address translation table. Losing this state would make it impossible to forward packets for existing connections.

Network model. We consider a deployment model where programmable switches are installed into the network fabric such that all traffic to be processed by an in-switch application traverses one of the programmable switches. This could be achieved in several different ways, depending on the network architecture. In a typical data center architecture (Figure 5.1), this could be achieved by using the switches on all core or all aggregation-layer switches.³ All traffic entering or leaving a cluster, for example, would traverse one of these switches. Alternatively, an operator might deploy a cluster of programmable switches as dedicated “NF accelerators”, explicitly routing traffic through them; this approach is similar to how software load balancers [90, 163] are deployed today.

State partitioning. We assume that application state is partitionable using some key derived from the packet header, and that each packet’s processing uses only state from the associated partition. In many cases, such as for the NAT example, the key will be the IP 5-tuple, and, hence, we use “partition” and “flow” interchangeably. However, other applications may use different partitioning, *e.g.*, partitioning on VLAN ID to detect heavy-hitter flows for a particular tenant.

We also assume that the network is configured to provide best-effort affinity such that packets from the same partition usually arrive at the same switch. Standard layer-3 routing protocols such as Equal-Cost Multi-Path routing (ECMP) provide this property when they are configured to use the partition key as their hash key.

Primer on programmable switches. Programmable switch architectures used today, *e.g.*, Intel Tofino [53], use a limited amount of on-chip memory (*e.g.*, SRAM and TCAM) to provide a variety of stateful object abstractions, including tables, registers, meters, and counters. Applications can use these to keep state across multiple packets, such as the address translation table in the NAT example above. In the ingress and egress match-action pipeline, objects are allocated in each stage and accessed by packets via ALUs. These objects are also accessible by the switch control plane

²Other applications maintain only soft state in the switch and provide their own failure recovery mechanisms. These are not the focus of our work, though RedPlane could perhaps help simplify their design or improve recovery performance.

³In principle, RedPlane could be deployed on top-of-rack (ToR) switches, but it is potentially less useful. If each rack has one ToR switch, and it fails, connectivity to the servers in that rack is lost. RedPlane can restore the switch state onto a different rack, but depending on the application that may not be useful. However, if there are two ToR switches per rack, RedPlane would be useful.

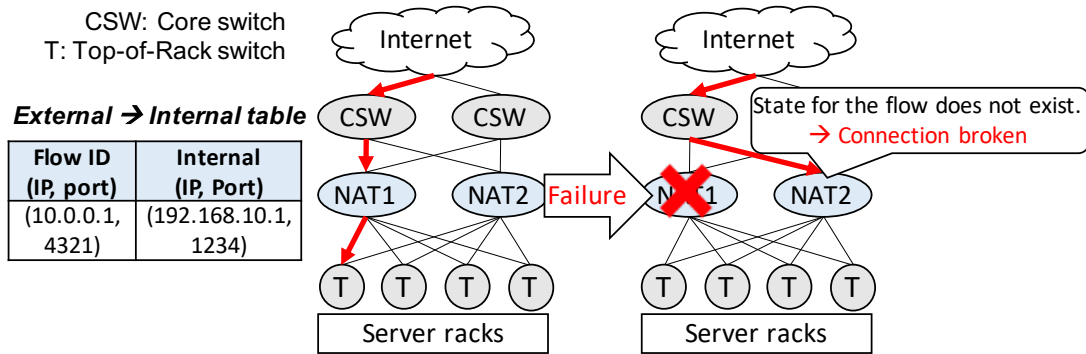


Figure 5.1: Impact of switch failures on in-switch NATs.

State access	Applications	Impact of switch failure
Read-centric	NAT	Connection broken
	Stateful firewall	Connection broken
	Load balancer [153]	Connection broken
	SYN flood defense [203]	Dropping valid packets
Write-centric	Super-spreader detection [186]	Inaccurate detection
	Heavy-flow detection [146]	Inaccurate detection
Mixed-read/write	SGW in EPC [177]	Active session broken
	In-network sequencer [137]	Incorrect sequencing
	Per-object routing [139, 206]	Choosing wrong servers
	In-network key-value store	Losing key-value pairs

Table 5.1: Examples of stateful in-switch applications and impact of switch failures.

through the ASIC-to-CPU PCIe channel which has a limited bandwidth ($O(10\text{ Gbps})$) compared to the ASIC’s per-port bandwidth ($O(100\text{ Gbps})$). In addition, the ASIC provides other built-in functionality such as packet replication, recirculation, and mirroring for more advanced packet processing. While we use Tofino-based programmable switches for our work, we believe our design can be implemented on other programmable switch ASICs since hardware capabilities leveraged in RedPlane’s switch data plane (*e.g.*, packet mirroring) are general features supported by most switch ASICs.

5.1.1 Impact of Switch Failures

Switches can fail, either by a switch failing entirely (a fail-stop model), or by individual links losing their connectivity. Measurement studies in production data centers have shown that such switch failures are prevalent. For example, in Microsoft’s data center, 29% of customer-impacting incidents are related to hardware failures including ASIC failure, fiber cuts, or power failures [143], and in Facebook’s data center, 26% of incidents are related to switch failures [152].

Switch failures can impact stateful applications in two ways. If a switch fails entirely, all application state it held is lost. Beyond that, a link failure or the failure of a *different* switch can impact many paths in the network [145], causing traffic to be rerouted [64, 102, 134]. Traffic that previously traversed one switch might be routed to a different one, where the appropriate state is unavailable. In the absence of this state, application processing can fail. For example, as illustrated in Figure 5.1, lacking the proper translation table entries, the NAT cannot forward packets for existing connections, breaking open connections *en masse*. Indeed, this is a serious problem – software-based stateful load balancers at cloud providers implement complex failover mechanisms [90, 163].

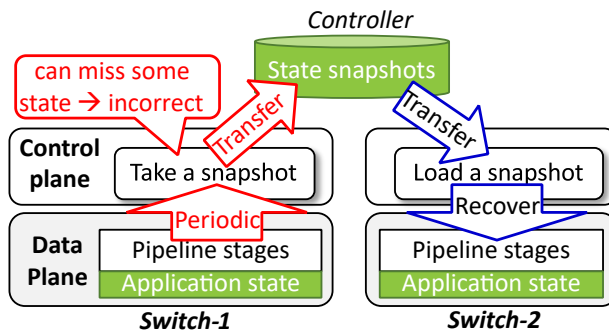
Beyond the conventional NFs (*e.g.*, NATs, load balancers, firewalls), there are several in-switch applications (shown in Table 5.1) that exhibit complex state access patterns. For example, many applications that are designed to enforce QoS policies (*e.g.*, rate limits) employ streaming algorithms (*e.g.*, sketching) to capture characteristics of traffic such as heavy-hitters [146, 186]. Switch failures lead them to make inaccurate decisions as the statistical data is lost. Such applications update state (*e.g.*, sketches) on every packet, so we call them *write-centric*. In contrast, many conventional NFs and DDoS defense systems (*e.g.*, SYN proxy) [197, 203] are *read-centric*.

Another group of applications have *mixed-read/write* state access patterns, typically with much less frequent updates than write-centric applications. One example in this category is NFs in the packet core for cellular networks (*e.g.*, Evolved Packet Core (EPC) for LTE) [37]. Packet core NFs such as a serving gateway (SGW) route users' data traffic from user devices to the Internet and vice versa based on per-user states (*e.g.*, forwarding state), which are updated when the control plane receives signaling messages (*e.g.*, device attached). To cope with the increasing volume of signaling traffic [7, 25],⁴ there have been recent efforts to accelerate the control plane functions by offloading them to the programmable data plane [11, 17, 164, 177]. For example, a SGW running on a switch maintains per-user tunnel endpoint IDs (TEIDs) to route packets, and this state is updated by signaling messages and read by data packets that are encapsulated with TEIDs. Thus, when a switch fails, since the SGW loses the state, it cannot forward packets for users, disrupting active connections. Affected users need to re-establish connections after the failure [60], increasing the service latency. Other applications that route requests in application-specific ways (*e.g.*, for databases [206] or key-value stores [139]) also fall into this category since they require state updates on every write (but not read) request.

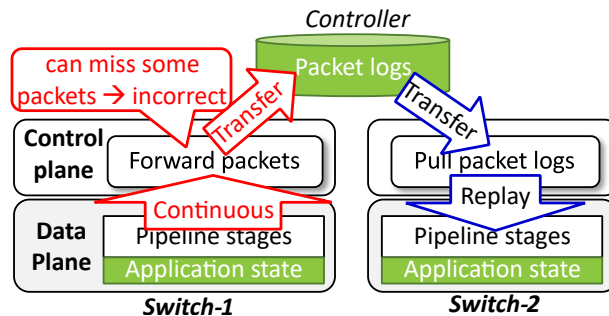
5.1.2 Existing Approaches and Limitations

We now examine classical fault tolerance mechanisms [98, 159, 191] and mechanisms tailored for network middleboxes [169, 179]. At a high level, these approaches can be categorized into three classes: (1) checkpoint-recovery, (2) rollback-recovery, and (3) state replication. All prior work targets server-based implementations. In what follows, we discuss why natural adaptations of these approaches to the switch environment fail to ensure correct behavior during failures.

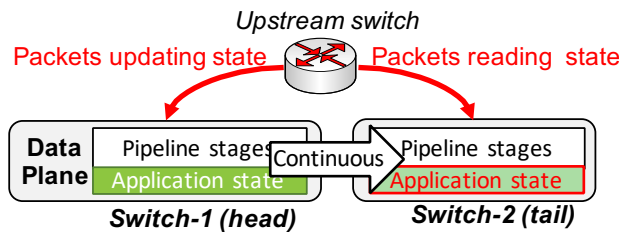
⁴Despite the growth, it is expected that signaling traffic rate is still much lower than that of data traffic (*e.g.*, 5% of data traffic [156]).



(a) Checkpoint-recovery: Switch control plane periodically snapshots and commits state to the controller. During this time, all packets must be buffered.



(b) Rollback-recovery: Each packet is forwarded to the control plane and logged by the controller.



(c) State replication: Switches replicate state to data plane memory using chain replication. Packets must be routed to the correct chain node. By design, the upstream switch at one-level below or above needs to route packets to a specific switch in a chain depending on the packet.

Figure 5.2: Highlighting why adapting existing approaches for fault tolerance fails for hardware switches.

Checkpoint-recovery. Checkpointing approaches periodically snapshot application state (e.g., an address translation table in NAT) and commit it to stable storage (e.g., [169]). When a failure occurs, the latest snapshot is populated on a backup node (i.e., an alternative switch in our context). Figure 5.2a illustrates a candidate implementation on switches using an external controller to store snapshots via the switch control plane. To achieve a consistent snapshot, data plane execution must be paused and packets buffered during the snapshot period. Limited data-to-control plane bandwidth in modern switch architectures makes this impractical.

Rollback-recovery. This approach, previously used for software middleboxes [179], logs every packet to stable storage and replays the traffic logs on a new device after failure to reconstruct application state. A natural implementation is sending every packet to the switch control plane, which logs it to the controller (Figure 5.2b). In principle, this approach can guarantee correctness if every packet is synchronously logged and replayed after a failure. However, the mismatch between the data traffic rate (Tbps) and the data-to-control plane bandwidth (Gbps) will result in many packets being dropped and will, thus, be incorrect.

State replication among switch data planes. Consider a state machine replication approach using chain replication [191], but applied to switch data planes (Figure 5.2c). Packets are forwarded through a sequence of switches, each of which updates its state and forwards the packet to the next switch in a chain. Only once the packet has reached the tail of the chain is it forwarded on its way to its destination. This is done entirely on the data plane, so it can function at high speed. This approach achieves correctness *only if* state updates are not lost. However, the state updates are delivered over an unreliable channel, and since the switch data plane cannot effectively support reliable transport protocols (*e.g.*, TCP) updates could be lost or reordered, violating correctness. Also, using one switch to replicate another switch’s state makes poor use of data plane-accessible switch memory – the most costly and limited resource. It also requires changes to the routing policy of the network since a packet needs to be explicitly routed to a specific switch in the chain depending on whether the packet updates state or not.

Takeaways. From the above discussion, we see two key takeaways. First, approaches that rely on the switch control plane must consider the mismatch between control and data plane speeds. Second, while switch data-plane-only approaches can provide good performance, they suffer three shortcomings: (a) incurring significant switch resource overhead; (b) making it difficult to reason about correctness due to unreliable communication channels between switch data planes; and (c) they may additionally constrain routing policies.

5.2 Overview

Our goal is to design a fault tolerance solution that provides the following four properties:

- **Correctness:** Switch failure should be transparent to applications: clients should not see state that would not be possible in the absence of a failure.
- **Performance:** Under failure-free operation, overhead for per-packet latency should be low (say, a few tens of μs).
- **Low resource overhead:** It should not consume switches’ limited compute and storage resources excessively.
- **Transparency to routing policies:** That is, we must allow a packet to update and/or read state regardless of the location of a switch where the packet is routed.

To this end, we present RedPlane, which provides an abstraction of fault-tolerant state storage for stateful in-switch applications. RedPlane provides an illusion of “one big fault-tolerant switch” – the behavior is indistinguishable from the same application running on a single switch

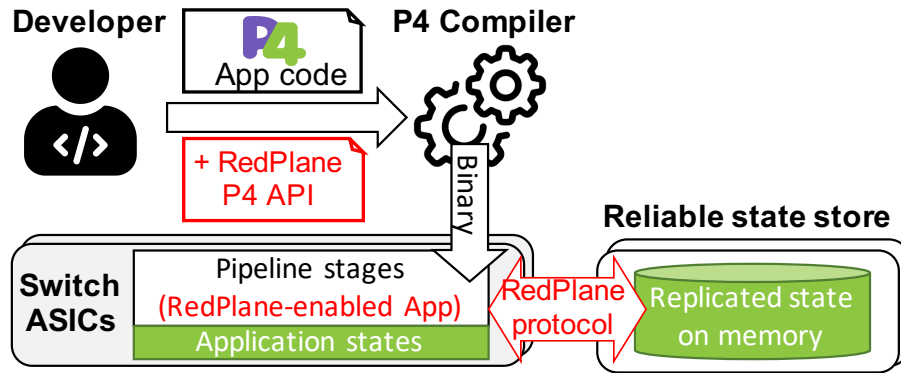


Figure 5.3: RedPlane overview highlighting extensions to traditional workflows for in-switch applications

that never fails. To achieve this, RedPlane continuously replicates state updates which can be restored without loss after a failure.

RedPlane takes a state replication approach with two defining characteristics: (1) the switch’s state replication mechanism is implemented entirely in the data plane, and (2) state storage is done through an *external state store*, a reliable replicated service made up of traditional servers. Property (1) means that the switch’s control plane is not required for state replication, avoiding the issues with the checkpointing and rollback-recovery approaches of [Section 5.1.2](#). Property (2) means that the replicated state is stored in commodity server DRAM, a relatively low cost storage medium compared to switch data plane memory. This avoids the high resource overhead of the state replication approach discussed in [Section 5.1.2](#).

While the idea of using servers’ memory as an external store is similar to recent work on TEA [129], it is important to note that TEA does not tackle fault tolerance. It focuses on the problem of resource augmentation to enable a switch to retrieve state stored in memory of servers. Furthermore, that design can only utilize servers *directly attached* to a top-of-rack switch. As such, their design does not tackle fault tolerance or provide provable correctness in the scenarios when multiple switches can access the store.

RedPlane provides a set of APIs ([Figure 5.3](#)) implemented in P4 [31], a language to specify data plane programs on programmable switches, to allow developers to easily integrate RedPlane with their stateful P4 applications. Once developers (re)write their applications using RedPlane APIs, the P4 compiler generates a binary of RedPlane-enabled applications loaded to the switch, which continuously replicates updates to the state store through the data plane.

Scope and limitations: In this work, our focus is on enabling fault tolerance for stateful applications with partitionable hard state, where a loss of state disrupts network or application functionality, shown in [Table 5.1](#). Applications only with non-partitionable state (*e.g.*, global counter) are beyond the scope of this work. Also, we assume that global state in an application (*e.g.*, a port pool in NAT) is sharded across and managed by state store servers. Other applications that need soft-state (*e.g.*, in-network caches or ML accelerators) do not require fault tolerance, but may benefit from RedPlane.

5.2.1 Design Challenges

While replicating state updates through the data plane to an external state store seems appealing, realizing this idea in practice presents some challenges:

- **Challenge 1. Providing correct replication in the data plane while tolerating unreliable communication** Traditional server-based replicated systems aim to provide strict correctness by ensuring not just linearizability but also that each operation is executed exactly once even in the presence of dropped or retransmitted messages [135, 142, 158]. To do so, they build on reliable communication channels like TCP. However, the switch data plane cannot support reliable communication, nor can it buffer significant amounts of traffic.
- **Challenge 2. Handling high traffic volume** Switch data plane operates at immense traffic volumes (up to a few billion packets per second [33, 50, 54]), in contrast to server-based systems handling a few million. If each packet that reads or updates state requires interacting with a server-based state store, the servers' capacity will rapidly be exceeded. It will also incur significant performance overhead.
- **Challenge 3. Being transparent to routing policies** A switch failure, recovery, or network routing change could cause traffic flows originally processed at a switch S1 to be routed to a different switch S2. However, since the routing decisions may be unpredictable, we cannot make assumptions on S2 or presuppose what backup routes will be taken. That is, we must be able to transparently migrate the relevant state from S1 to S2 irrespective of the location of S2. For instance, we need to make the NAT table entry available when packets for a particular connection are processed by a different instance.

5.2.2 Key Ideas

To tackle these challenges, we build on four key ideas:

- **Idea 1: Practical correctness for switch state (Section 5.3).** We define two correctness models based on the requirements of in-switch applications. The first, a strict consistency mode, is based on linearizability [111]. Because we observe that network applications are already designed to tolerate packet loss, we explicitly adopt the standard definition of linearizability, which permits operations that do not complete while still providing strong consistency for those that do. Second, since many write-centric applications (*e.g.*, monitoring using sketches [84]) accept approximate results, we propose a relaxed consistency mode that allows some state to be lost after a failure, but bounds the inconsistency.
- **Idea 2: Piggybacking output packets (Section 5.4.1).** Instead of buffering output packets using limited switch resources, we use the network itself as temporary storage by piggybacking packet contents on coordination messages.
- **Idea 3: Lightweight sequencing and retransmission (Section 5.4.2).** To cope with the unreliable communication channel between the switch data plane and the state store with low resource overhead, we employ a sequencing mechanism for protocol messages and devise a lightweight switch-side retransmission mechanism by repurposing the switch ASIC's packet mirroring feature.

- **Idea 4: Lease-based state ownership (Section 5.4.3).** To reduce the frequency with which the switch must coordinate with the state store, especially for applications with read-centric and mixed-read/write workloads, we adopt a lease-based mechanism inspired by prior work [100, 142, 158]. This allows us to avoid coordination with the state store for packets which need to read but do not modify state. At the same time, we ensure that all state updates are durably recorded before any of their effects are externalized, guaranteeing linearizability. This mechanism also serves as the means by which state is migrated between switches to support the transparency.

Taken together, these high-level ideas address the aforementioned challenges. First, the linearizability-based consistency model coupled with the piggybacking and lightweight sequencing and retransmission mechanism allows to replicate state reliably and correctly (Challenge 1). Second, the relaxed consistency and lease-based state ownership help cope with high traffic volume (Challenge 2). Lastly, the lease-based state ownership makes RedPlane transparent to routing policies (Challenge 3).

5.3 Correctness Model

RedPlane provides two levels of consistency, which applications can choose between based on their requirements. A *linearizable* mode provides strict guarantees, making the system indistinguishable from a single fault-tolerant switch. Because this has a high overhead for write-centric applications due to frequent coordination with the state store, RedPlane also offers a *bounded-inconsistency* mode that permits some state updates to be lost on switch failure, but guarantees a consistent view of switch state.

5.3.1 Preliminaries

By default, RedPlane provides *linearizability* [111], a correctness condition for concurrent systems. We model a stateful in-switch program as a state machine, where the output and next state are determined entirely by the input and current state:

Definition 1 (Stateful in-switch program). A stateful program P is defined by a transition function $(I, S) \rightarrow (O^*, S')$ that takes an input packet and the current state, and produces zero, one, or multiple output packets, along with a new state.

To simplify the definitions below, we will assume that each input packet p produces exactly one output packet $P(p)$; it is straightforward to extend them to the zero- or many-output case. This implies that the program’s behavior is determined entirely by the sequence of input packets, and in particular that it is deterministic and that packets are processed atomically. Although switch architectures are pipelined designs that process multiple packets concurrently [71], their compilers assign state to pipeline stages in a way that makes packet processing appear atomic [66].

The gold standard for replicated state machine semantics is single-system linearizability [111]. That is, that the observed execution matches a sequential execution of the program that respects the order of non-overlapping operations. To adapt linearizability for in-switch programs, we first redefine a history in terms of packet processing:

Definition 2 (History). A history is an ordered sequence of events. These can be either input events I_p , in which a packet p is received at a RedPlane switch, or output events O_p in which the corresponding output packet is output by a RedPlane switch.

Note that it is possible for there to be input events I_p without the corresponding output O_p , if the processing of p is still in process or due to a failure. We discuss this in depth next.

5.3.2 Linearizable Mode

Definition 3 (Linearizability for a stateful in-switch program). A history H is a linearizable execution of a program P if there is a reordering S of the input events in H such that (1) the value for each output event O_p that exists in H is given by running P on the input events in S in sequence, and (2) if O_x precedes I_y in H then I_x precedes I_y in S .

Here, S is the apparent sequential order of execution.

Linearizability is fundamentally a safety property, not a liveness one: it specifies what output values are acceptable, but does not guarantee that all operations complete. It is possible for a packet to be received and (1) update switch state, but produce no output, or (2) neither update switch state nor produce any output. Definition 3 reflects this: a packet with an input event but no output event can still appear in the sequential order S . If it precedes the processing of other packets, then they see the effects of its state update. If it appears at the end of S , it has no visible effect on system state.

While these anomalies comport with the definition of linearizability, most replicated systems aim to provide a stronger property: that every operation is executed exactly once and returns its result to the client. Ensuring this requires several protocol-level mechanisms: typically, clients retry requests that do not receive a response, and replicas keep state to detect duplicate requests and resend the responses without executing them twice [135, 142, 158]. As we see (Section 5.4.2), these techniques are not feasible in our environment.

Accordingly, RedPlane takes a different approach: it *explicitly permits* these two types of anomalies. While this may seem surprising, it matches the semantics of modern networks. The two cases correspond to a packet being lost (1) between the RedPlane switch and its destination or (2) between the source and the RedPlane switch, respectively. Network applications must already tolerate lossy networks, so they are resilient to such losses.

Relaxing the definition of correctness enables a tractable implementation. By not requiring the system to achieve complete reliability, our protocol may drop packets during failover, or if messages between a switch and the state store are lost. In these scenarios, an input packet or its output may be lost. Of course, dropping too many packets is undesirable for performance reasons; such loss events are rare.

5.3.3 Per-flow Linearizability

In most in-switch programs, some or all state is associated with a particular flow – a subset of traffic identified by a unique key, *e.g.*, an IP 5-tuple, VLAN ID, or an application-specific object ID. For example, each translation table entry in a NAT is tied to a specific flow based on an IP 5-tuple. For many applications, per-flow state is the only state that needs to be consistent or

fault tolerant – either because there is no global state, or because global state can tolerate weaker consistency, *e.g.*, traffic statistic counters that need not be precise. RedPlane generally provides consistency for per-flow state (consistency for global state is optional):

Definition 4 (Per-flow linearizability). A history H is per-flow linearizable if, for each flow f , the subhistory H_f for the packets in flow f is linearizable.

As long as programs use only per-flow state, per-flow linearizability is *the same* as global linearizability, because linearizability is a *local* (*i.e.*, composable) property [111]. The benefit of operating on a per-flow level is that it means synchronization between states associated with different flows are not required. As we show in Section 5.4, this allows RedPlane to distribute execution of a program across multiple switches: each has the state associated with certain flows, and can process packets for those flows. This matches the way many applications are deployed in practice, *e.g.*, a NAT will be deployed to a cluster of switches, using ECMP for load balancing. Because this load balancing is done on a per-flow granularity, each switch is responsible for performing translation for a subset of flows, and does not need access to the translation table for the other flows.

5.3.4 Bounded-inconsistency mode

RedPlane’s linearizable mode uses a synchronous replication protocol (Section 5.4.1), which can induce high overhead for write-centric applications. However, we observe that many write-centric applications in programmable switches operate in contexts where approximate results are acceptable, *e.g.*, monitoring using sketches [84] or Bloom filters [70]. For these applications, RedPlane offers a *bounded-inconsistency* mode that has lower overhead.

In this mode, RedPlane takes periodic snapshots of data plane state and replicates them asynchronously. This means that upon switch failure, the most recent state updates can be lost. However, RedPlane ensures that the system recovers to a consistent state from within a time interval ϵ . RedPlane’s consistent snapshot mechanism ensures that the state after recovery reflects an actual state of the system, which simplifies reasoning about the correctness of complex data structures.⁵ In Section 5.4.4, we describe how we address key challenges in realizing this mode in RedPlane.

5.4 RedPlane Design

Now, we describe the RedPlane protocol that realizes our linearizable and bounded-inconsistency modes. We begin with an overview of the protocol and explain how we address practical challenges.

⁵Although the bounded-inconsistency mode may affect properties of some approximate data structures (*e.g.*, no false negatives in Bloom filters), since it bounds the inconsistency within ϵ , developers or network operators can easily reason about the potential inconsistency.

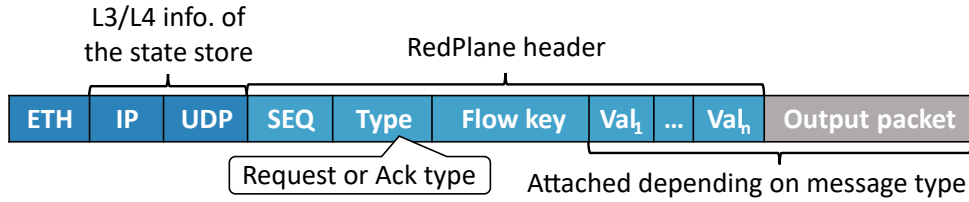


Figure 5.4: RedPlane state replication protocol packet format.

5.4.1 Basic Design

As shown in Figure 5.3, RedPlane consists of (1) an external state store built on commodity servers and (2) a RedPlane-enabled application running on the switch data plane. In this section, we describe how the components work together via the state replication protocol.

For clarity of exposition, we start with simplifying assumptions: there is no packet loss or reordering between switches and the state store, switches do not fail while messages are in transit, and packets for a flow are routed to only one switch at a time. We revisit these assumptions in Section 5.4.2 and Section 5.4.3.

External state store:

The external state store is an in-memory key-value storage system. We partition it across multiple shards by flow – identified by an IP 5-tuple or other key. Each state store shard can be replicated using conventional mechanisms and we do not seek to innovate here as many existing key-value stores meet our needs (e.g., [48, 141, 160]). Specifically, our prototype is a simple in-memory storage server implemented in C++ that uses chain replication [191] with a group size of 3.

Basic replication protocol:

A RedPlane-enabled application replicates state updates to the state store by exchanging protocol messages formatted as shown in Figure 5.4. It uses standard UDP and IP headers to address messages to the state store or the switch using their respective IP addresses. The RedPlane header consists of a sequence number, a message type, and a flow key. Depending on the message type, it can also include flow state and an output packet. We will discuss these fields shortly. Note that we assign an IP address to each RedPlane switch and use it for routing requests and response packets between state store servers and RedPlane switches. This works with general L3 routing protocols including ECMP and BGP.

As an illustrative example to help understand the protocol, we consider a per-flow counter application shown in Figure 5.5. This application updates or reads the state for each packet. In the example, there are two switches and a state store. We have multiple packets in each flow f , with the n^{th} packet denoted as pkt_n^f . This example illustrates a case where the Switch-1 initially handles f , but after its failure, the flow is rerouted to the Switch-2.

State initialization or migration (Step ① or ④ in Figure 5.5). When the application receives a packet that belongs to a flow it has never seen before (e.g., pkt_1^{f1}), it needs to send a *state*

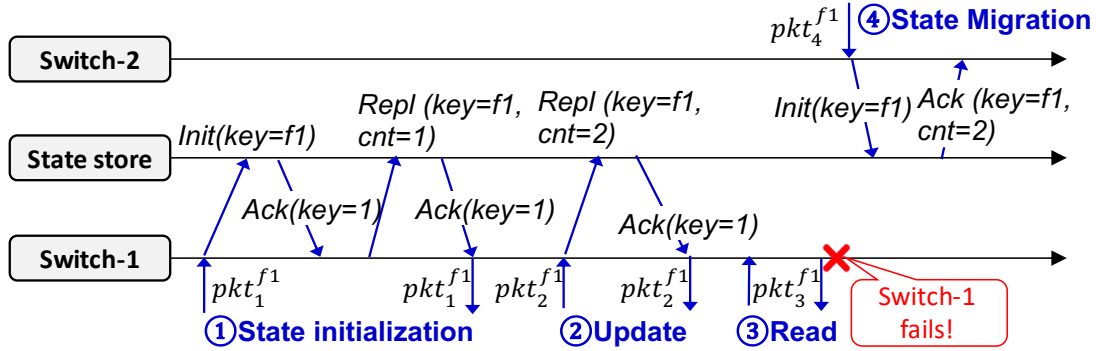


Figure 5.5: Basic workflow of RedPlane state replication protocol. “Repl” indicates a state replication request. pkt_n^f indicates n^{th} packet of a flow f .

initialization request. It identifies the corresponding state store server by hashing the flow key (e.g., IP 5-tuple), and looking up the corresponding server IP and UDP port from a preconfigured table.

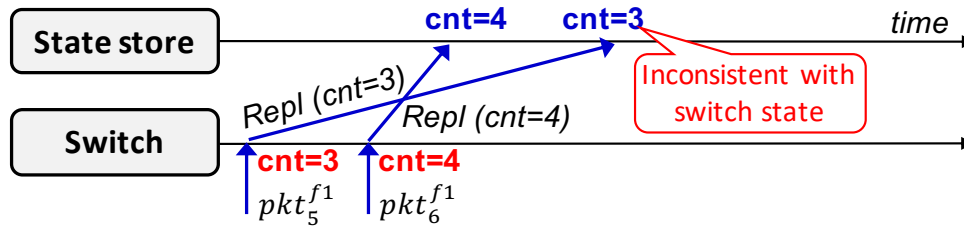
There are two possible cases: (1) the flow is new and so has no state, or (2) the flow state previously existed on a failed switch, and a packets for that flow are now being routed to a switch on an alternative path (*i.e.*, failover). In case (1), upon receiving the request, the state store initializes its storage for the state and sends a response back to the switch (Step ①). In case (2), since the state store already has the flow state, it sends a response containing the latest state (Step ④).

Upon receiving the response, the application installs the returned state into the corresponding switch memory. For stateful memory registers, this can be done entirely in the data plane. On the Tofino architecture, updates to match tables or certain other resources need to be done through the switch control plane. In this case, RedPlane routes the processing through the control plane. This can introduce additional latency (we measure this in Section 5.6.1). However, many in-switch applications already require a control plane operation on a new flow (e.g., to install a new translation mapping in a NAT), in which case the added overhead is minimal.

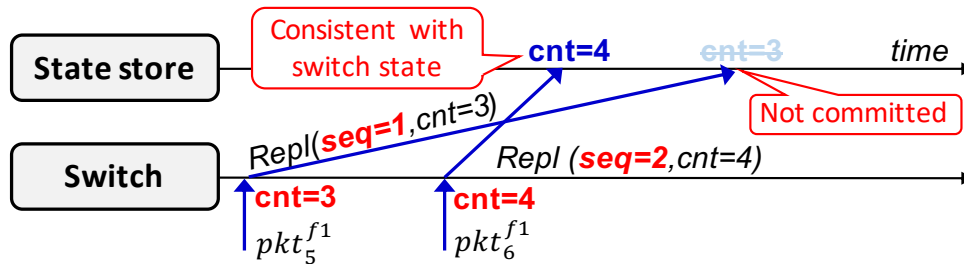
Reading or updating state (Step ② or ③ in Figure 5.5). Once the state has been initialized, the application can read the state value (*i.e.*, the counter in our example) directly (Step ③). When it updates the state (*i.e.*, the counter value), RedPlane sends a *replication request* with the new value to the state store. This message is generated entirely through the data plane. The state store applies the update, and sends a *replication reply* message (Step ②).

Piggybacking output packets. When the application updates the state, RedPlane should not allow an output packet to be released until the state has been recorded at the state store – otherwise, the update could be lost during a switch failure, violating correctness. This requires the output packet to be buffered until the replication reply is received.

Unfortunately, the switch data plane does not have sufficient memory to buffer packets in this way (and various other constraints on how memory can be accessed make it unsuitable for storing complete packet contents). RedPlane instead piggybacks the packet onto its replication request message, and the state store returns it in its reply. When the reply is received, RedPlane



(a) Out-of-order delivery of requests causes inconsistency.



(b) Request sequencing serializes replication requests.

Figure 5.6: Serializing out-of-order requests with sequencing. Counter values (cnt) in red and blue indicate the state on the switch and the state store, respectively.

decapsulates and releases the packet. In effect, this uses the network and the memory on the state store as a form of delay line memory – trading off network bandwidth, which is plentiful on a switch, for data plane memory, which is scarce.

Note that it is possible to receive packets that read state when there are in-flight replication requests for the state. In this case, the packets are buffered in the same way through the network (with a special RedPlane request type) until a switch receives a response for the latest replication request.

While our basic design provides correctness under the simplified assumptions, we find that in more realistic environments, it may not be able to guarantee correct behavior. In the following sections, we describe potential challenges, and how we extend the basic design to address them.

5.4.2 Sequencing and Retransmission

To guarantee correctness, replication requests must be *successfully* delivered and replicated *in order* at the state store. For example, the replication request (**Step ②** in Figure 5.5) must be delivered in order. However, successful in-order delivery is not guaranteed in a best-effort network between switches and the state store.

Figure 5.6a illustrates why such unreliability in the network can be problematic. We use the same per-flow counter as an example. Each time the counter is incremented, RedPlane sends the new value to the state store. If the state store just processes updates in the order they are received,

a reordering could cause a later counter value to be replaced with an earlier one. Request loss can cause a similar issue.

A traditional replication system, like chain replication, might address this by relying on a reliable transport protocol like TCP. Unfortunately, it is not practical to implement a full TCP stack on the switch data plane – if it is possible at all, it would excessively consume data plane resources.

Our approach. Instead of implementing a full-fledged reliable transport on a switch data plane, we choose to build a simpler UDP-based transport with mechanisms that deal with possible packet reordering and loss. First, to handle out-of-order state replication request messages, we employ a mechanism called *request sequencing* [137], which assigns a per-flow monotonically increasing sequence number to each request message. The state store uses this sequence number to avoid applying updates out of order (Figure 5.6b).

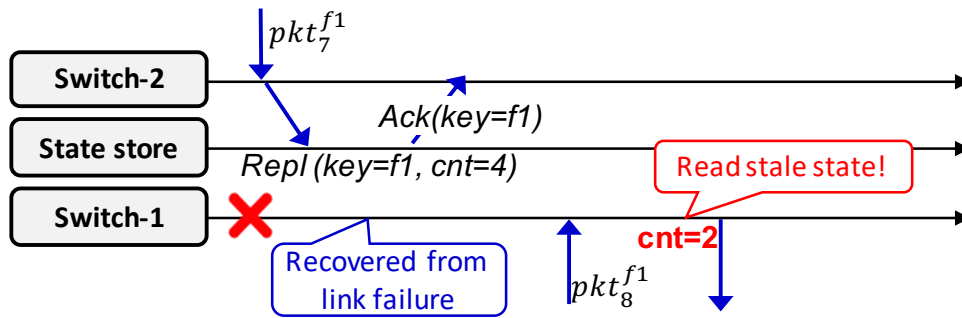
Second, to cope with lost replication requests or responses, we develop a mechanism for *request buffering*. RedPlane buffers replication requests and retransmits them if it does not receive a reply before a timeout. We implement this by repurposing the egress-to-egress packet mirroring capability of switch ASICs. When RedPlane sends a replication request, it mirrors a copy with the current timestamp as metadata. When the mirrored request enters the egress pipeline and it has not been acknowledged by a response with the same or a higher sequence number, RedPlane checks whether the request has timed out by comparing the current timestamp to the timestamp in its metadata. If it has timed out, it resends the request to the state store. Otherwise, it mirrors the request again without ending the request to the state store.

As discussed previously, buffering a full packet payload is challenging on a switch due to memory limitations. Instead, RedPlane buffers *only state updates* (i.e., the RedPlane header) – not the piggybacked output packet by truncating the packet. This reduces the amount of data that needs to be mirrored. A consequence of this is that if a replication request or its response is dropped, the output packet will be lost. This is permitted by our linearizability-based correctness model: it is indistinguishable from the output packet being sent and dropped in the network. The state updates must be retransmitted, however, because subsequent packets processed by the switch may see the new version, and thus it must be durably recorded. We measure the overhead of request buffering in Section 5.6.4.

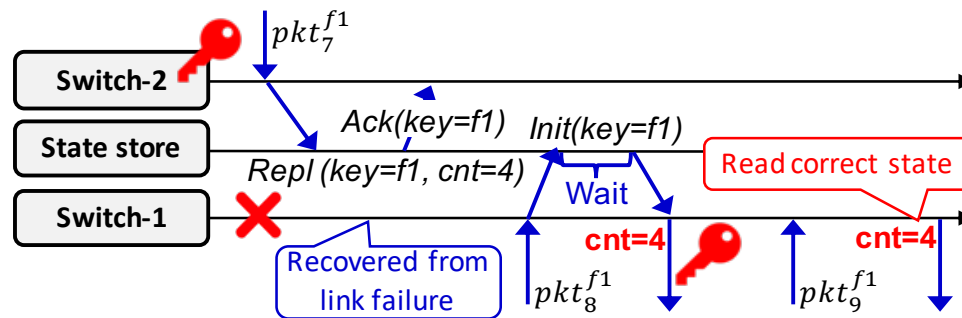
5.4.3 Lease-based State Ownership

What if multiple switches attempt to process packets for a particular flow at the same time, especially during failover or recovery? The protocol in Section 5.4.2 will not be correct in this case, when there are concurrent accesses to the same state. Figure 5.7a illustrates why. After Switch-1 has a link failure (but does not lose its state, which is $cnt=2$), packets are routed to an alternate, Switch-2. If Switch-1 recovers, a packet may read its old state, a violation of linearizability.

Our approach. RedPlane ensures that only one switch can process packets for a given flow at a time using *leases*, a classic mechanism for managing cached data in file systems [100] and replicated systems [142, 158]. Figure 5.7b illustrates this. If a packet wants to access state, but



(a) Stale state access after a switch recovers from link failure.



(b) Only one switch holds a lease (red key) on state at a time.

Figure 5.7: Consistent state access for multiple switches.

the state is not available at the switch, it first requests a lease for the flow. The state store grants a lease for a specific time period (1 second in our prototype) only if no other switch holds an active lease on the same flow state. The lease time is renewed each time the switch sends a replication request for that flow; switches that frequently read but infrequently update state can send explicit lease renewal requests. Our prototype does so every 0.5 seconds.

5.4.4 Periodic Snapshot Replication

As described in Section 5.3, RedPlane offers bounded-inconsistency mode for write-centric applications that permit approximate results, *e.g.*, monitoring using sketches [186] or Bloom filters [197]. In this section, we describe how we realize it in the switch data plane.

For such applications, RedPlane replicates *snapshots* of state *asynchronously* and *periodically*. Every T_{snap} seconds, a snapshot of the current state is sent to the state store, while output packets are released without waiting for replication to complete.

However, realizing this approach entirely in the data plane is challenging. While data structures often consist of multiple entries (*e.g.*, slots in sketches), the switch is architected, and the P4 language is designed, to allow access to a single entry per register array per packet. Also, building hardware that could atomically copy entire register arrays would be costly.

To address this challenge, we employ a *lazy snapshotting* approach. We maintain two copies of the data structure that are lazily synchronized with each other. These are interleaved in the switch’s register arrays so that each array index contains two entries, one from each copy. Two metadata registers are used to indicate which entry at each index is the *active* copy. The first, a 1-bit flag, is toggled when a snapshot is taken. The second, a 1-bit register array, represents whether that index has been updated since the current snapshot started.

To take a snapshot, we flip the flag and read values from the now-inactive copy. Meanwhile, when packets arrive and update the array, one of two operations occur. The first packet to update an index synchronizes the two copies and then updates the active copy. Later packets simply update the active copy. This allows us to take a consistent snapshot of the entire structure while incoming packets continue to update it. Additional snapshots must wait for the current one to complete. We describe the pseudocode of our mechanism in [Appendix B](#).

Replication is achieved using the switch ASIC’s packet generator. We configure it to generate a batch of packets every T_{snap} seconds. To replicate a data structure with n entries, we generate a batch of n packets, each with a unique ID p_i . The ID in each packet is used to address the i th entry in the data structure and copy its value into a RedPlane replication protocol header. Note that while RedPlane asynchronously replicates snapshots, it still guarantees successful replication with its sequencing and retransmission mechanisms.

5.4.5 Protocol Correctness

RedPlane’s replication protocol provides per-flow linearizability defined in [Section 5.3](#). Due to space constraints, we give only a brief sketch of the reasoning here. The lease protocol ensures that at most one switch is executing a program for a particular flow at a time. The sequencing, retransmission, and buffering protocol ensure that an output packet is never sent unless the corresponding state update has been recorded and acknowledged by the state store.

During non-failure periods, RedPlane provides per-flow linearizability because the single switch processing packets for a flow operates linearizably, but some output packets may be lost (due to dropped replication traffic with piggybacked messages). After a failover, the new switch receives a state version at least as new as the most recent output packet from the old switch. This satisfies the linearizability requirement that any packet sent after these output packets were observed follow it in the apparent serial order of execution. We also wrote a TLA+ specification of the linearizable mode to model-check the above property ([Appendix C](#)).

Our periodic snapshot replication guarantees that the system recovers to a consistent state from within a time bound ϵ (*i.e.*, bounded inconsistency) by tracking the time since the last successful replication; if the time bound is exceeded, an application-specific action may be taken (*e.g.*, dropping further packets or treating the switch as failed).

5.5 Implementation

Our prototype implementation is available in our repository [[59](#)].

Data plane. We implement RedPlane’s data plane components in P4-16 [31] (≈ 1192 lines of code) and expose them as a library of P4 control blocks [31, §13], which form the RedPlane APIs that developers can use to make application state fault tolerant. We compile RedPlane-enabled applications to the Intel Tofino ASIC [53] with P4 Studio 9.1.1 [45]. We implement key functions such as lease request generation, lease management, sequence number generation, and request timeout management, using a series of match-action tables and register arrays. We evaluate the additional resource usage in Section 5.6.4. As mentioned in Section 5.4.2, we implement request buffering via the mirroring and truncation capabilities of the switch ASIC, which allows us to buffer only the replication protocol data and discard the original payload. We implement a basic sketch that supports lazy snapshotting; developers can modify it to implement similar data structures such as Bloom filters.

Control plane. We implement the switch control plane in Python and C++. Its main function is to initialize and migrate (if available) state for the data plane by processing corresponding responses forwarded by the data plane component.

State Store. Our contribution is in the fault tolerance protocol design and switch components. As such, our state store prototype is built based on readily available libraries and simple implementations. We implement RedPlane’s state store in C++ for Linux servers. It uses Mellanox’s kernel-bypass raw packet interface [5] for optimized I/O performance. To ensure reliability in the presence of server failures, we implement chain replication [191] using a group of 3 servers located in different racks.

Applications. To demonstrate the applicability of RedPlane, we implement various applications in P4 described below. Figure 5.8 illustrates how the P4 implementation of RedPlane-enabled NAT looks like. Developers need to include the P4 file of RedPlane core APIs (line 1) and the P4 file of their original application code (line 2). Lines highlighted in red shows initialization and the use of the RedPlane ingress and egress control block instances (line 5, 9, 20, and 24). And the lines highlighted in bold blue indicates modules of the original NAT program (line 6 and 11). Since NAT does not update state in the data plane (i.e., read-centric), no modification is needed to their original P4 implementation. Other applications described below can be implemented in a similar way.

(1) *NAT*: The NAT implementation uses RedPlane to implement a fault-tolerant per-5-tuple address translation table and available port pool. Since the port pool is a shared by different flows, it is sharded across state store servers and managed by them. The state is updated when a TCP connection is established from an internal network.

(2) *Firewall*: The stateful firewall adds fault-tolerance to a per-5-tuple TCP connection state table using RedPlane. Its state is updated when a TCP connection is established from an internal network.

(3) *Load balancer*: The load balancer maintains a per-5-tuple server mapping table; we make it fault-tolerant using RedPlane. It also uses a server IP pool, which is shared state. When a new TCP connection is established from an external network, the state is updated.

(4) *EPC-SGW*: We also implement a simplified serving gateway (SGW) used in cellular networks, a mixed-read/write application. It maintains per-user tunnel endpoint ID state. The state is updated by signaling messages and read by data packets.

```

1  #include "redplane_core.p4" // RedPlane core API
2  #include "nat.p4"          // developer's NAT program
3
4  control Ingress (headers hdr, metadata meta) {
5      RedPlaneIngress() redplane_ingress;
6      NAT_Ingress() nat_ingress;
7      L3_Routing() l3_routing;
8      apply {
9          redplane_ingress.apply(hdr, meta);
10         if (meta.is_normal_pkt == true) {
11             nat_ingress.apply (hdr, meta);
12         }
13         if (meta.is_normal_pkt == true ||
14             meta.is_piggybacked == true) {
15             l3_routing_ingress.apply (hdr, meta);
16         }
17     }
18 }
19 control Egress (headers hdr, metadata meta) {
20     RedPlaneEgress() redplane_egress;
21     apply {
22         if (meta.is_redplane_req == true ||
23             meta.is_redplane_ack == true) {
24             redplane_egress.apply(hdr, meta);
25         }
26     }
27 }
28
29 Pipeline(
30     IngressParser(),
31     Ingress(),
32     IngressDeparser(),
33     EgressParser(),
34     Egress(),
35     EgressDeparser()
36 ) pipe;
37
38 Switch(pipe) main;

```

Figure 5.8: The main part of P4 implementation of RedPlane-enabled NAT.

(5) *Heavy-hitter (HH) detection*: We implement a heavy-hitter detector using count-min sketches [84] as an example of write-centric applications; there are 3 sketches, each consisting of 64×32 -bit slots indexed by a hash of the IP 5-tuple. We implement separate sketches per VLAN ID, assuming that the network operator wants to enforce different policies for each cloud tenant. Since sketches are an approximate data structure which can be replicated asynchronously, we use periodic snapshot replication.

(6) *Per-flow counter*: To demonstrate RedPlane's worst-case performance, this application counts packets forwarded for each IP 5-tuple. State is updated for every packet and synchronous replication must be used.

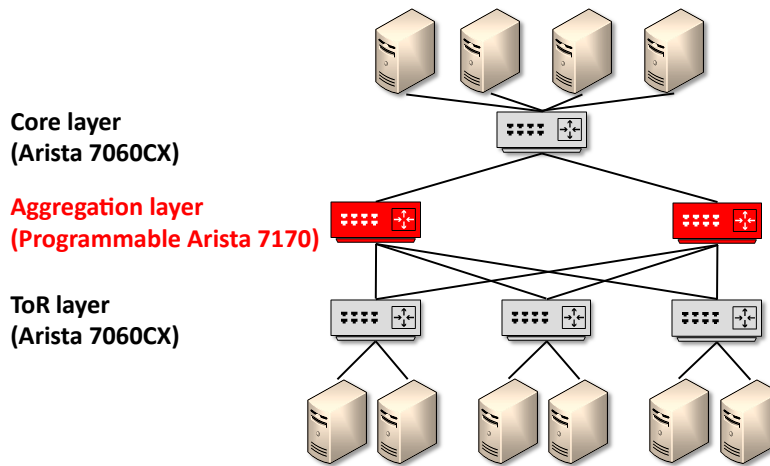


Figure 5.9: Three-layer network testbed for experiments.

5.6 Evaluation

We evaluate RedPlane on a testbed consisting of six commodity switches (including two programmable ones) and servers using both real data center network packet traces and synthetic packet traces. Our key findings are:

- In failure-free operation, RedPlane adds no per-packet latency overhead for applications that are read-centric or replicate state asynchronously. For write-centric applications in linearizable mode, RedPlane incurs $8 \mu\text{s}$ per-packet overhead (Section 5.6.1).
- In failure-free operation, the throughput of read-centric applications is not degraded. For write-centric applications, the throughput is bottlenecked by state store performance in linearizable mode, but periodic snapshot replication reduces the overhead. Similarly, RedPlane incurs almost no bandwidth overhead for read-centric applications and small overhead for write-centric in bounded-inconsistency mode even at scale (Section 5.6.2).
- After a switch failure, RedPlane-enabled applications access their correct state and recover end-to-end TCP throughput within a second (Section 5.6.3).
- RedPlane provides these benefits with little resource overhead as it consumes $<14\%$ of ASIC resources (Section 5.6.4).

Testbed setup. We build on a three-layer network testbed consisting of six commodity switches (including two programmable ones) and servers, as shown in Figure 5.9. The aggregation layer has two 64-port Arista 7170 Tofino-based programmable switches [33] running stateful applications written in P4. The core and ToR switches run 5-tuple-based ECMP routing to route packets to end hosts even when one aggregation switch fails. Each ToR switch has two servers connected, and four additional servers attached to the core switch emulate hosts outside the datacenter. The state store runs on one server in each rack. All servers are equipped with an Intel Xeon Silver 4114 CPU (40 logical cores), 48 GB DRAM, and a 100 Gbps Mellanox ConnectX-5 NIC, running Ubuntu 18.04 (kernel version 4.15.0). We repeat each experiment 100 times unless otherwise noted.

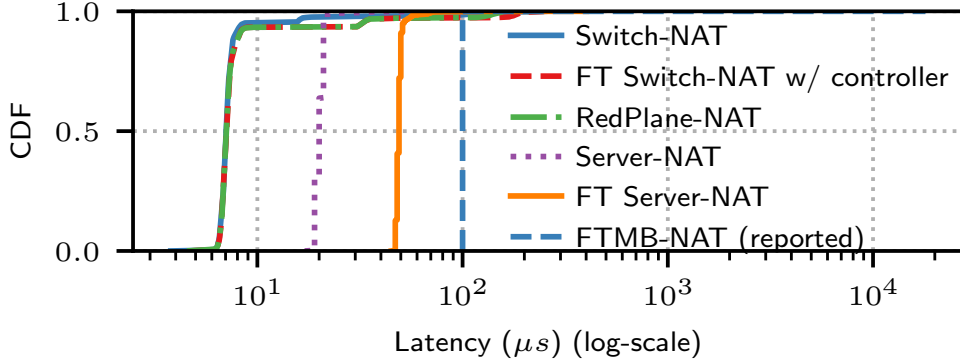


Figure 5.10: End-to-end RTT when RedPlane-NAT processes packets *vs.* other approaches.

5.6.1 Latency in Normal Operation

First, we evaluate the per-packet latency overhead introduced by RedPlane under failure-free operation for the 5 applications in Section 5.5. To measure the processing latency, we have each application send packets back to a sender node and track the RTT of each packet. We replay publicly available packet traces from a real data center and enterprise network [2, 3] to generate 100,000 packets and measure the processing latency of each packet. The packet sizes vary (64–1500 bytes) in the real traces. To evaluate EPC-SGW, we inject a signaling packet for every 17 data packets, following statistics used in previous studies [156, 167].

Overhead of RedPlane. As an exemplar application, we evaluate the per-packet latency for a NAT in RedPlane⁶ and compare it with baseline implementations: (1) NAT written in P4 without fault-tolerance (Switch-NAT), (2) NAT written in P4 with controller based fault-tolerance (Switch-NAT w/ an external controller)⁷ (3) NAT implemented on a CPU server without fault-tolerance (Server-NAT), (4) NAT implemented on a server with fault-tolerance (FT Server-NAT), and (5) FTMB-NAT which uses rollback-recovery for server-based middleboxes [179].⁸ For switch-NAT w/ controller, RedPlane-NAT, and server-NAT, we enable chain replication for the controller, state store, and NAT instances, respectively.

Figure 5.10 shows the CDF of the per-packet latency distribution. Compared to Switch-NAT, which is expected to have the lowest latency, RedPlane-NAT shows the same 50th and 90th percentile latency ($7 \mu s$ and $8 \mu s$, respectively), meaning that there is no overhead. This is because for NATs, packets except for the first packet of each flow only require state (*i.e.*, address translation table) to be read. Both Switch-NAT and RedPlane-NAT show a high 99th percentile latency ($110 \mu s$ and $142 \mu s$, respectively), mainly due to the overhead introduced by our control plane implementation; in Switch-NAT, the first packet of every flow is forwarded to the switch control plane to create and insert a new entry to the translation table. RedPlane-NAT has additional overhead since it needs to request a lease from the state store before updating

⁶We choose NAT to compare results with those reported in prior work [179].

⁷We implement a simple external controller to emulate SDN controller-based approaches (*e.g.*, Morpheus [175] and Ravana [121]), which communicates with the switch control plane via a 1 Gbps management channel.

⁸We use the latency reported in the original FTMB aper [179] since we were not able to get its full implementation.

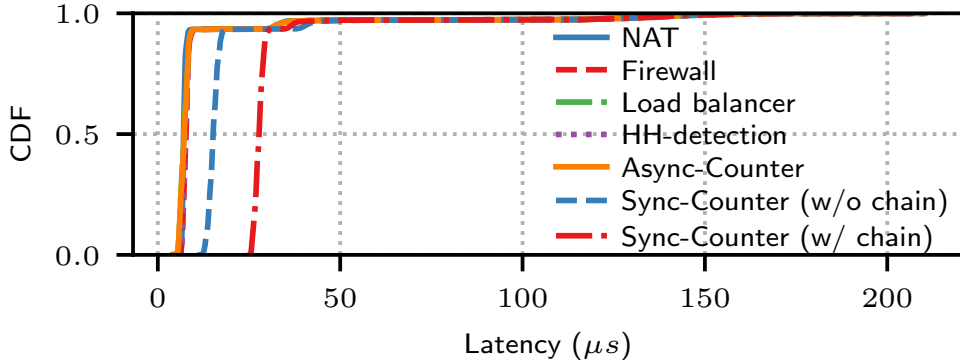


Figure 5.11: End-to-end RTT for RedPlane-enabled applications. All applications have chain replication enabled for the state store. For Sync-Counter, we also show its overhead without chain replication.

state. Switch-NAT with the external controller incurs higher 99th percentile latency ($185 \mu s$) due to the communication overhead between the switch control plane and the controller and between controller instances (for chain replication) over the slower management network. Server-based versions (FT Server-NAT and FTMB-NAT) have $7\text{--}14\times$ higher median latency compared to the switch-based approaches, as packets need to traverse additional hops in the network and they have inherent performance limitations.

Impact on different applications. Next, we evaluate the per-packet processing latency overhead of different RedPlane-enabled applications. As shown in [Figure 5.11](#), RedPlane-enabled NAT, firewall, load balancer, EPC-SGW, and heavy-hitter (HH) detection, all have the same $8 \mu s$ median latency, identical to that without fault-tolerance. The NAT, firewall, and load balancer are read-centric and update state only when a new flow is created; EPC-SGW is mixed-read/write, and updates state on signaling packets whose frequency is 5% of data packets. HH detection, although it is write-centric, performs periodic state replication asynchronously, so it does not affect the latency. On the other hand, since Sync-Counter updates state and replicates updates *synchronously* for every packet, it adds an additional latency of $20 \mu s$ to every packet. $12 \mu s$ of this overhead is due to the 3-way chain replication used to tolerate state store server failures.

5.6.2 Bandwidth Overheads

To evaluate network bandwidth overheads, we inject 64-byte packets from three traffic generation servers at ≈ 207.6 Mpps⁹ which is the maximum rate that our traffic generator can achieve.

Additional bandwidth consumed. In this experiment, we instrument each application to count the number of bytes it sends and receives, including both original packets and protocol message packets. [Figure 5.12](#) shows the ratio of bandwidth used for RedPlane messages to the total traffic. For read-centric applications including NAT, firewall, load balancer, we see that there is almost no bandwidth overhead since RedPlane generates protocol messages only for the first packet of

⁹Each server generates packets at ≈ 69.2 Mpps.

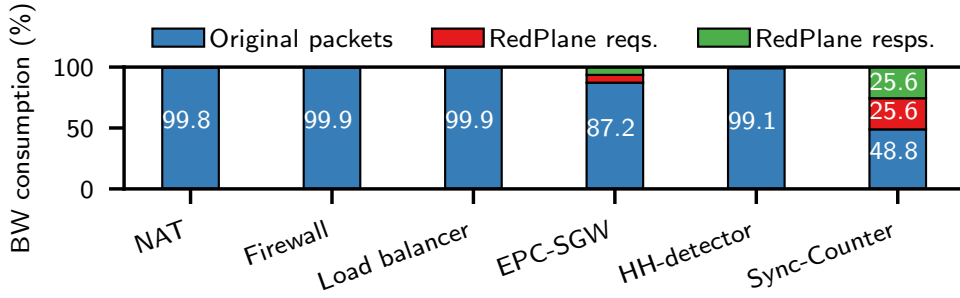


Figure 5.12: RedPlane replication bandwidth overhead.

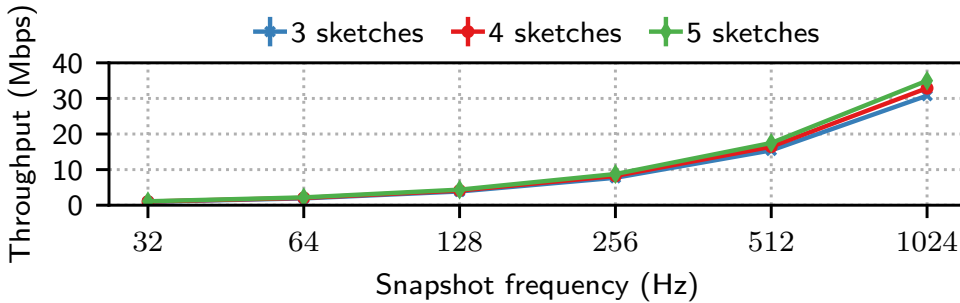


Figure 5.13: Impact of the frequency of snapshotting on RedPlane-enabled HH-detector.

each flow. For EPC-SGW, RedPlane incurs 12.8% overhead since it generates protocol messages for signaling packets, and some of data packets are buffered through the network as described in Section 5.4.1. For HH-detector, which asynchronously replicates a snapshot of state for every 1 ms, RedPlane incurs negligible overhead. We also measure the absolute bandwidth overhead for different snapshot frequencies and number of sketches as shown in Figure 5.13. For a 1 ms period, it consumes 34.16 Mbps (13.8%). Even with 5 sketches, this is lower than the bandwidth overhead for Sync-Counter (51.2%) because in the latter case RedPlane requests and responses contain both headers and original payload. This result implies that in an extreme case where an application replicates state updates synchronously for every packet, achieving fault-tolerance is expensive. We also analyze the bandwidth overhead at scale (*i.e.*, a topology with more RedPlane switches) for all 6 applications using our analytical model-based simulation, and the result is consistent with Figure 5.12 in terms of the percentage overhead.

Throughput impact on applications. In this experiment, we measure the throughput of RedPlane-enabled applications and compare it with the same applications without fault tolerance. We send 64-byte packets from three servers, one from each rack, to one of servers attached to the core switch at ≈ 207.6 Mpps. In our testbed, the link between an aggregation and a core switch becomes the bottleneck, and we observe that the maximum forwarding rate the aggregation switch can achieve is around 122.5 Mpps. Figure 5.14 shows the median throughput of each application with and without RedPlane. Obviously, applications achieve the maximum throughput without RedPlane. With RedPlane, read-centric (NAT, firewall, and load balancer) applications and applications that replicate state updates asynchronously (HH-detector) can

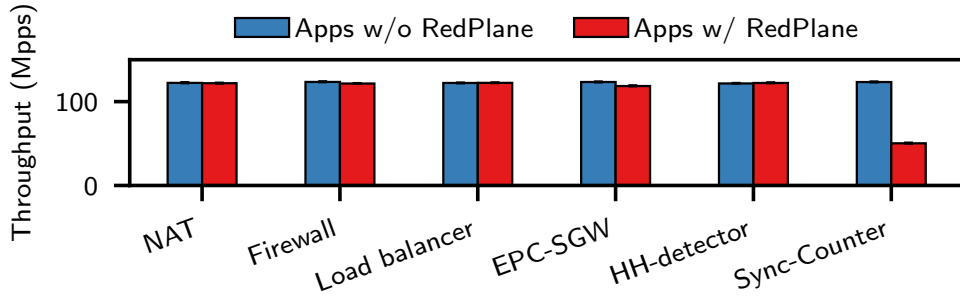


Figure 5.14: Impact of RedPlane on data plane throughput of RedPlane-enabled applications.

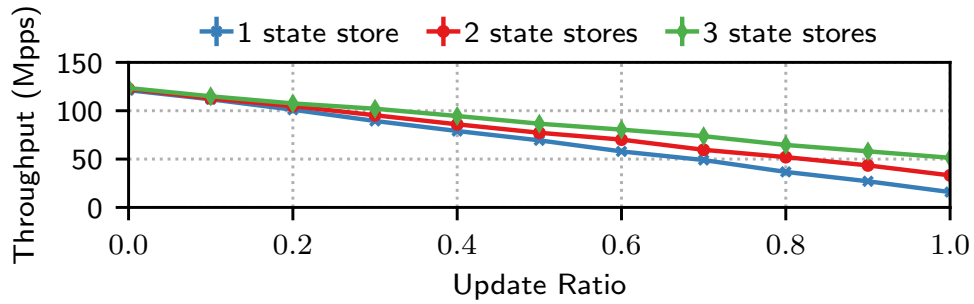


Figure 5.15: Impact of update ratio on data plane throughput of the RedPlane-enabled key-value store.

achieve the same throughput as their non fault-tolerant counterparts. The RedPlane-enabled EPC-SGW achieves a slightly lower throughput than that of its counterpart, mainly due to some data packets buffered through the network during the replication. The throughput of Sync-Counter becomes nearly half that of its counterpart: we find that it is bottlenecked by the performance of the state store. This suggests that applying a strict consistency mode degrades the throughput of write-centric applications as they are also affected by the performance of the state store.

Varying update ratios. While most of existing in-switch applications are read-centric or perform asynchronous replication, incurring little overhead, it is important to understand the maximum throughput of applications characterized by different read/write (*i.e.*, update) ratios. For this experiment, we write a simple in-switch key-value store in P4 with RedPlane and generate packets consisting of custom header fields that indicate an operation (read or update), a key, and a value (for updates). We use the same setup as the previous experiment and let each server generate packets based on a predefined update ratio with uniformly distributed random keys. Figure 5.15 shows that as the update ratio increases, the throughput degradation depends on the number of state store servers; by adding more servers, we can achieve higher throughput.

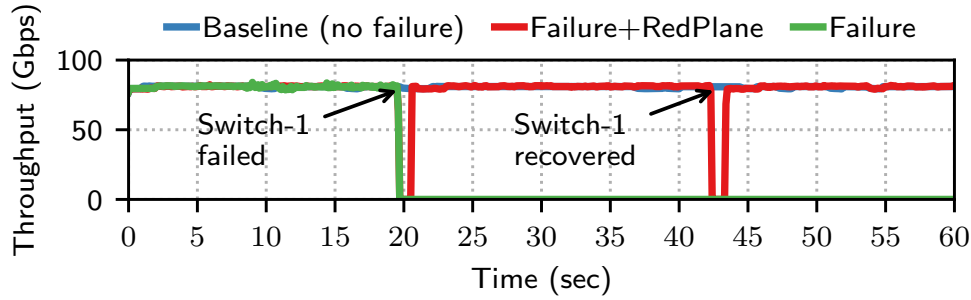


Figure 5.16: End-to-end throughput changes during failover and recovery with and without RedPlane.

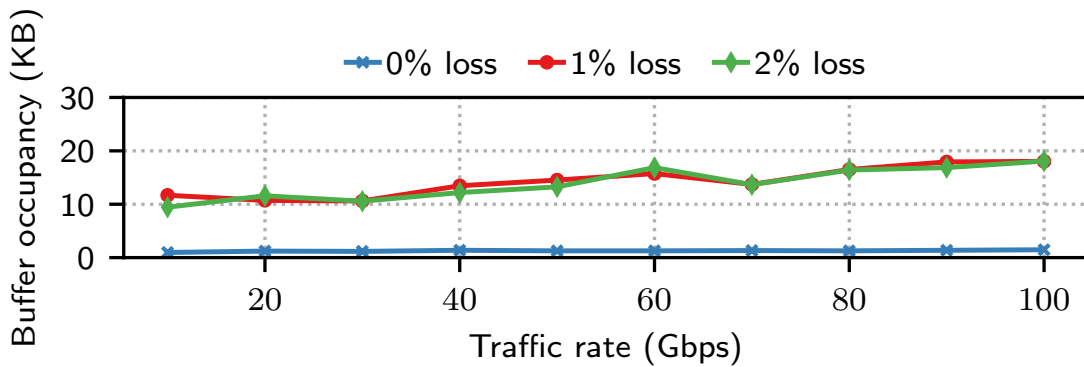


Figure 5.17: Switch packet buffer occupancy due to request buffering.

5.6.3 Failover and Recovery

Next, we measure how fast the end-to-end performance can be recovered by RedPlane in the presence of switch failure and recovery. We run `iperf` [40] to measure between two servers, attached to a core switch and a ToR switch respectively. All traffic passes through a NAT running on the programmable switches. We compare changes in TCP throughput when (1) there is no failure (Baseline), (2) one switch fails without RedPlane (Failure), (3) one switch fails while using RedPlane (Failure+RedPlane).

Figure 5.16 shows the results. In a network without RedPlane, when Switch-1 fails, packets are rerouted to another switch and dropped, breaking the TCP connections. In contrast, RedPlane-enabled NAT successfully maintains high throughput when the switch fails and recovers after short disruptions (0.9 and 1.0 seconds). This recovery time is affected both by the core switch’s failure detection/rerouting time and RedPlane’s lease period (set to 1 second here). Control plane and state store optimizations could further reduce this.

5.6.4 RedPlane Switch ASIC Resource Usage

Packet buffer usage. In this experiment, we evaluate the overhead of our request buffering mechanism (Section 5.4.2). Since RedPlane buffers a replication request until it receives a reply

Resource	Additional usage
Match Crossbar	5.3%
Meter ALU	8.3%
Gateway	9.9%
SRAM	13.2%
TCAM	11.8%
VLIW Instruction	5.5%
Hash Bits	3.7%

Table 5.2: Switch ASIC resources used by RedPlane.

corresponding to the request from the state store, it consumes some amount of the switch packet buffer. Since there is no precise way of measuring the buffer usage in real-time, we instead use the queue depth information provided by the switch ASIC to estimate the upper bound of the buffer occupancy.¹⁰ Specifically, we assume a write-centric application where every incoming packet issues a request (*i.e.*, the most demanding scenario). And we let each request packet record its queue depth information to a P4 register in the data plane and read it from the control plane for every second and take the maximum value. We generate packets from a traffic generation server while varying the traffic rate and the request loss rate.¹¹ Figure 5.17 shows the result. When there is no request loss, the buffer occupancy is less than 1.5 KB even at 100 Gbps traffic rate. As we increase the request loss rate, the buffer usage also grows; when the traffic rate is 100 Gbps and $\approx 2\%$ of requests are lost, our buffering mechanism consumes at most 18 KB, which is acceptable for a given a few tens of MB of the packet buffer in the switch ASIC.

Table 5.2 shows the additional switch ASIC resource consumption of RedPlane for 100K concurrent flows (using the P4 compiler’s output), expressed relative to each application’s baseline usage. Overall, there are ample resources remaining to implement other functions along with RedPlane.¹² RedPlane uses TCAM to implement acknowledgment processing and request timeout management, which need range matches. In terms of scale vs. number of concurrent flows, only the SRAM usage would increase proportional to the number of flows as it stores per-flow information (lease expiration time, current sequence number, and last acknowledged sequence number).

5.7 Related Work

Fault-tolerance for in-switch applications. Recent efforts have shown that offloading to programmable switches enhances performance. For example, offloading the sequencer [137], key-value cache [115, 147], and coordination service [116] improves the performance of distributed

¹⁰It is a per-packet queue depth measured when a packet is dequeued from the buffer, and the Tofino ASIC provides this information as an intrinsic metadata that can be accessed at the egress pipeline.

¹¹We emulate the request loss by dropping requests at a certain probability at the switch.

¹²Match Crossbars are used for implementing the ‘matching’ part of match-action tables. Meter ALUs perform stateful operations on registers. Gateways perform ‘if-else’ conditions in the control flow.

systems. However, these applications can lose their state due to switch failures. RedPlane can help make them fault-tolerant or simplify their designs.

Fault-tolerance and state management for server-based NFs. Fault-tolerance for server-based NFs or middleboxes has been addressed by prior systems like Pico [169] and FTMB [179]. When an NF instance fails, the state of the failed NF is recovered through checkpoint or rollback recovery on a new NF instance. These approaches cannot be applied directly to the switch data plane (Section 5.1.2). Previous work on state management for stateful NFs uses local or remote storage to manage NF state [97, 170, 196]. However, these APIs target planned state migration rather than unplanned failures. Similar work (again, targeting planned migration) has also been proposed for router migration [124].

External memory for switches. TEA (Chapter 3) shares RedPlane’s approach of using servers’ memory as external storage for switch state [129], but towards a different goal: allow switches to handle state larger than their on-device memory. It does not address fault tolerance or multi-writer consistency.

Switch-based reliability protocols. Other recent work runs coordination protocols between switches to build reliable storage [86, 116]. Our goal is conceptually different – to replicate state for in-switch applications rather than provide a networked storage service – but uses some similar mechanisms, like network sequencing [137].

5.8 Summary

While many recent efforts have demonstrated the potential benefits of running datacenter functions on programmable switches, we argued that there is one critical missing piece in current designs, which is fault tolerance. To address this issue, in this chapter, we presented RedPlane, which provides a fault tolerant state store abstraction for in-switch applications. We formally defined a linearizability-based correctness model for a replicated switch data plane state and build a practical replication protocol based on it. Our evaluation with various stateful applications on a real testbed showed that RedPlane can support fault-tolerance with minimal performance and resource overheads and enable end-to-end performance to quickly recover from switch failures.

Chapter 6

Conclusions

In this thesis, we have argued that by designing the right abstractions and runtime environments for programmable data plane devices, we can build an in-network computing platform that supports resource elasticity and fault resiliency without any hardware modifications.

With **Table Extension Architecture (TEA)** ([Chapter 3](#)), we explored a new approach that leverages DRAM on remote servers available in a typical data center, with a focus on network functions (NFs) as an application. TEA provides a virtual table abstraction that allows NFs on programmable switches to look up large virtual tables built on external DRAM. Our approach enables switch ASICs to access external DRAM purely in the data plane without involving CPUs on servers and the switch control plane. TEA provides low and predictable table lookup latency scalable throughput with additional servers (*e.g.*, 138 million lookups per second with 8 servers).

With **ExoPlane** ([Chapter 4](#)), we argued that an on-rack switch resource augmentation architecture that augments a programmable switch with other programmable network hardware, such as smart NICs, on the same rack could be an affordable and incrementally scalable solution. To realize this vision, we designed and implemented ExoPlane, an operating system for on-rack switch resource augmentation to support multiple concurrent stateful applications. ExoPlane provides in-network applications with low and predictable latency, scalable throughput, and fast failover while achieving these with small resource overheads and no or little modifications on applications.

With **RedPlane** ([Chapter 5](#)), we designed and implemented a fault-tolerant state store for stateful in-network applications. RedPlane provides in-switch applications consistent access to their state, even if the switch they run on fails or traffic is rerouted to an alternative switch. We address key challenges in devising a practical, provably correct replication protocol and implementing it in the switch data plane. RedPlane incurs negligible overhead and enables end-to-end applications to recover from switch failures rapidly.

Before concluding, we summarize a few of the lessons learned during the course of this thesis work and sketch future research directions in this space.

6.1 Putting All The Pieces Together

Since the goal of our three systems (TEA, ExoPlane, and RedPlane) is to provide resource elasticity or fault resiliency for stateful in-network applications, the natural question is: *Can the three systems be used simultaneously for a stateful in-network application?* Although this thesis has not shown an integrated use case, one could build a memory-intensive and fault-tolerant in-network application by using the P4 interfaces provided by each system as they are independent. For example, in our ongoing work, we leverage the ExoPlane and RedPlane APIs to build a defense system that can mitigate volumetric multi-vector attacks (*e.g.*, distributed denial of service (DDoS) attacks). In this example system, each attack mitigation function running on programmable switches is a stateful program, and ExoPlane and RedPlane make multiple concurrent functions fault-tolerant using resources on external smart NICs and commodity servers. One challenge in using different systems' APIs simultaneously is that the restriction and lack of modularity in the P4 language require developers to manually combine two API implementations at the P4-code level, which is an error-prone and tedious task. As future work, we plan to address this challenge with the compiler or language support to automate the integration process (Section 6.3).

6.2 Lessons Learned

Throughout this thesis research, we have explored various programmable networking hardware devices and built our runtime environments on top of them. In this section, we share some of the lessons we have learned and discuss what we need for future-proof in-network computing.

Need for runtime reconfigurability for data plane devices. For most of the programmable devices we have used, we find that an ability to reconfigure the functionality of the data plane (*e.g.*, replacing a program A with a program B or partially reconfiguring a program's logic) is essential, but it is one of the missing features in today's data plane devices. For example, from our experience, for Intel Tofino-based programmable switches [53] and Netronome Agilio CX smart NICs [30], there is a few 10s seconds of "downtime" when loading a new binary (or firmware) to the ASIC or NPUs. Even worse, FPGA-based NICs such as Intel N3000 [41] takes a few tens of minutes to load a new bitstream! During this reconfiguration period, devices cannot serve any incoming traffic, significantly limiting the availability of in-network computing platforms, and we had to consider this constraint when designing our systems carefully. For example, in ExoPlane, we assumed that a set of applications or traffic workload distributions do not change on the small timescale (*e.g.*, for every minute) so that the ExoPlane planner does not need to reconfigure the switch and external devices frequently. While such an assumption is reasonable in our setting, we admit that it can limit other possible deployment models (*e.g.*, traffic workload distributions change for every minute). To make in-network computing platform future-proof, we believe that future hardware device design should consider runtime reconfigurability into account. One possible way is to partially reconfigure the ASIC or NPUs so that while reconfiguring the main logic, the remaining "active" portion of the device can process traffic as a fallback mode.

Need for better control-to-data plane communication interfaces. While our primary focus in this thesis is on in-network applications running in the network data plane, we realize that

the interface between the control to the data plane, which is used to manipulate objects in the data plane (*e.g.*, tables and register arrays) from the control plane, can be critical for overall performance. As we analyzed in RedPlane, many applications involve at least one control-to-data plane operation to initialize or update the data plane objects, which takes a few tens of *ms* in the case of Intel Tofino-based programmable switches and Netronome Agilio CX NICs. Also, in terms of throughput, we find that it supports up to a few thousand operations per second. The root cause of this is the delay from the PCIe interface between the control plane and the data plane and the control plane’s OS kernel driver stack. In most cases, these control plane operations are rare (*e.g.*, for the first packet of each flow), so we often assume that this delay becomes a one-time cost. However, we observe that in some cases, these delays can be non-negligible and impact application performance. For example, if some adversarial traffic generates new flow at a high rate, the control plane cannot keep up with the request rate. Also, if one needs to implement a fast feedback loop between the control and data plane, a few tens of *ms* of latency would not be acceptable. Thus, we need a better and faster communication interface between the control and the data plane to support such scenarios.

Need for better understanding of proprietary blackbox compilers. While implementing the data plane components of our systems and various applications running on them, we often face compile-time errors from vendor-provided proprietary backend compilers (*e.g.*, Intel Tofino P4 compiler), which complain about failing to allocate resources or violate some rules without providing detailed reasons even when there seem to be available resources. In such cases, since we do not know the exact reasons why the compiler fails, we have to manually fine-tune our P4 implementations to try different resource allocations, which is tedious and could be error-prone. We believe that open-sourcing the compiler or even part of it (*e.g.*, heuristic resource allocation algorithm it uses) would be helpful for building and debugging in-network applications more systematically.

6.3 Future Directions

Analysis of adversarial workloads. As network devices become programmable, their data and control plane functionality can be more intricate than conventional fixed-function devices, potentially leading to more potential security vulnerabilities. Attackers can exploit such vulnerabilities to bombard the data (*e.g.*, impacting false-positive rates of probabilistic data structures [83]) or control planes (*e.g.*, forwarding an excessive volume of traffic to the device control plane) by generating adversarial traffic patterns that trigger abnormal behaviors. Thus, TEA-, ExoPlane-, RedPlane-enabled applications can also be vulnerable to adversarial workloads. Unfortunately, there is no systematic method for identifying and fixing such possible vulnerabilities in in-network applications. One potential approach is automatically generating adversarial traffic workloads by statically analyzing the data and control plane implementation of in-network applications that can run on heterogeneous devices.

Manageable in-Network computing at scale. While in-network computing can bring us many benefits, a fundamental question remains: how can we deploy and manage in-network applica-

tions in a large-scale network without affecting the rest of the network? Answering this question is critical for large-scale deployments of in-network applications. This question involves many challenges, including concurrent updates of multiple device data planes, deployment of multiple concurrent applications, state management, and device heterogeneity. Addressing these challenges is critical for large-scale deployments.

Automatically enabling elasticity and resiliency with language and compiler supports. TEA, ExoPlane, and RedPlane provide application developers with APIs in the form of P4 control blocks so that they can use our elastic and resilient resource abstractions when building their applications. However, due to the constraints of the current P4 language (*e.g.*, program codes cannot be fully modularized as in other high-level programming languages such as C++) and vendor-provided proprietary compilers, our systems currently require some manual efforts from developers. For example, in TEA and RedPlane, developers have to manually insert the data plane components of our runtime environments and combine them with original application logic manually. Ideally, this process should be automated via support from a language and a compiler. For example, to make a match-action table fault-tolerant, developers annotate the table object, and a compiler front-end detects the annotation and automatically combines our runtime logic with the application logic would be one possible direction. Moreover, while our implementation is mainly on the data plane, we find that integrating the control plane component of our runtime environments with an application’s control plane component requires some engineering efforts. Thus, automatically combining the control plane logic while not affecting the correctness of each of them is an important problem to tackle.

Verifying the performance and correctness of in-network applications. One of the missing pieces in today’s in-network computing is the ability to verify the performance and correctness of applications running on different types of multiple devices in the network at scale. As we add more functionalities to individual network devices, it can affect the performance or correctness of other applications. While there have been several works on verifying the correctness of network configurations and data plane programs, they focus on the correctness of either the network’s control plane or a single program running on a single data plane device. Also, they do not consider the performance aspect of data plane programs. Thus, we need a way to systematically verify the correctness and performance of multiple applications simultaneously running on multiple devices in the network. Such a verification tool can even be useful for testing our approach at scale. Conceptually, our runtime environments add a logical “plane” to the network that allows accessing external resources at runtime while physically sharing network resources with the existing data plane (*e.g.*, network links). While we have demonstrated the feasibility of our approaches and evaluated their performance in our current deployment model, the performance and correctness verification tools will be essential to extend it at larger scale.

Designing next-generation switch architectures as future-proof in-network computing platforms. This thesis focuses on applications that require capabilities supported by the current programmable switching ASICs (*e.g.*, not requiring packet payload inspection). Looking forward, we see that there are increasing demands from applications that could benefit from in-network computing but cannot be realized on today’s switches. For example, applications that require

inspection of packet payloads, such as intrusion detection over fragmented packet streams (*e.g.*, TCP streams), or middlebox functions over encrypted packet streams, are not possible on switches due to the lack of an ability to process packet payloads. While other devices in the network (*e.g.*, a software switch [20, 36] and FPGA-based smart NICs) could implement some of the necessary functionality, the design of modern switches makes it challenging to incorporate this external functionality seamlessly. Furthermore, we observe that the current inflexible design of the switch data plane makes it difficult to keep up with these evolving application demands. Traditionally, switching ASICs are designed to support only stateless or simple stateful functionality such as packet forwarding and access control in a synchronous manner at line-rate (*e.g.*, a few Tbps). Because of this, even today, they are equipped with a limited amount of fast and expensive memory (*e.g.*, SRAM and TCAM) and compute (*e.g.*, simple ALUs and hashing units) on-chip resources, which are complex or infeasible to be extended once the chip is manufactured. Due to these constraints, they cannot support evolving applications. Thus, we need to rethink how the switch architecture as an in-network computing platform is designed to support evolving application demands.

Appendix A

Simplified Codes of TEA-enabled NF Implementations

```
1  #include "tea_core.p4"
2  ...
3  control Ingress (headers hdr, metadata meta) {
4  ...
5  LookupHandler() lookup_handler;
6  ServerResolver() server_resolver;
7  MemResolver() mem_resolver;
8  apply {
9  ...
10     lookup_handler.apply(hdr, meta);
11     if (is_ext.apply().hit) { //packet from external?
12         if (meta.lookup_md.found == true) {
13             forward.apply();
14         } else {
15             if (meta.lookup_md.remote_miss == false) {
16                 server_resolver.apply(meta);
17                 mem_resolver.apply(meta);
18             }
19         }
20     }
21     ...
22 }
23 }
24 control Egress (headers hdr, metadata meta) {
25     LookupRequestor() lookup_req;
26     apply {
27         if (meta.lookup_md.found == false &&
28             meta.lookup_md.remote_miss == false) {
29             lookup_req.apply(hdr, meta);
30         }
31     }
32 }
```

Figure A.1: Firewall.

```
1  #include "tea_core.p4"
2  ...
3  control Ingress (headers hdr, metadata meta) {
4  ...
5  action update_server_addr () {
6      hdr.ipv4.dstIP = meta.lookup_md.serverIP;
7  }
8  ...
9  LookupHandler() lookup_handler;
10 ServerResolver() server_resolver;
11 MemResolver() mem_resolver;
12 apply {
13     ...
14     lookup_handler.apply(hdr, meta);
15     if (meta.lookup_md.found == true) {
16         update_server_addr();
17         forward.apply();
18     } else {
19         server_resolver.apply(meta);
20         mem_resolver.apply(meta);
21     }
22     ...
23 }
24 }
25 control Egress (headers hdr, metadata meta) {
26     LookupRequestor() lookup_req;
27     apply {
28         if (meta.lookup_md.found == false) {
29             lookup_req.apply(hdr, meta);
30         }
31     }
32 }
```

Figure A.2: Load balancer.

```

1  #include "tea_core.p4"
2  ...
3  control Ingress (headers hdr, metadata meta) {
4    ...
5    action nat_ext_to_int () {
6      hdr.ipv4.dstIP = meta.lookup_md.pip;
7      hdr.ipv4.dstPort = meta.lookup_md.pport;
8    }
9    ...
10   table nat {
11     key = {
12       meta.lookup_md.dir: exact;
13     }
14     actions = {
15       nat_ext_to_int;
16       nat_int_to_ext;
17       drop;
18     }
19   }
20   ...
21   LookupHandler() lookup_handler;
22   ServerResolver() server_resolver;
23   MemResolver() mem_resolver;
24   apply {
25     ...
26     lookup_handler.apply(hdr, meta);
27     if (meta.lookup_md.found == true) {
28       nat.apply();
29       forward.apply();
30     } else {
31       server_resolver.apply(meta);
32       mem_resolver.apply(meta);
33     }
34     ...
35   }
36 }
37 control Egress (headers hdr, metadata meta) {
38   LookupRequestor() lookup_req;
39   apply {
40     if (meta.lookup_md.found == false) {
41       lookup_req.apply(hdr, meta);
42     }
43   }
44 }

```

Figure A.3: Network address translator.

```

1  #include "tea_core.p4"
2  ...
3  control Ingress (headers hdr, metadata meta) {
4    ...
5    action update_server_addr () {
6      hdr.ipv4.dstIP = meta.lookup_md.serverIP;
7    }
8    ...
9    LookupHandler() lookup_handler;
10   ServerResolver() server_resolver;
11   MemResolver() mem_resolver;
12   apply {
13     ...
14     lookup_handler.apply(hdr, meta);
15     if (meta.lookup_md.found == true) {
16       update_server_addr();
17       forward.apply();
18     } else {
19       server_resolver.apply(meta);
20       mem_resolver.apply(meta);
21     }
22     ...
23   }
24 }
25 control Egress (headers hdr, metadata meta) {
26   LookupRequestor() lookup_req;
27   apply {
28     if (meta.lookup_md.found == false) {
29       lookup_req.apply(hdr, meta);
30     }
31   }
32 }

```

Figure A.4: VPN gateway.

Appendix B

Lazy Snapshotting Algorithm in RedPlane

[Algorithm 5](#) shows the pseudocode for lazy snapshotting described in [Section 5.4.4](#). We implement this logic in P4 to provide a basic sketch with 64×32 -bit slots. As explained in [Section 5.5](#), we implement count-min sketches using three of this sketch.

Algorithm 5: Lazy snapshotting

```
/* 1-bit variable indicating the current active buffer */
1 active_buffer ← 0;
/* array of 1-bit variables indicating which buffer has has lastly been updated for a
   certain slot */
2 last_updated_buffer[0..REGISTER_SIZE] ← 0;
/* two copies of the replicated data structure (e.g., a sketch in this example) */
3 pair<int,int> sketch [0..REGISTER_SIZE] ← 0;
4 Upon receiving a packet (pkt):
/* is this the first pkt of a snapshot read burst? */
5 if pkt.type = SNAPSHOT_READ and pkt.index = 0 then
/* if so, swap the active buffer */
6 | active_buffer ← swap_active_buffer();
7 else
/* if not, get the current active buffer */
8 | active_buffer ← get_active_buffer();
/* which buffer was lastly updated for this index? */
9 last_updated_buffer_for_index ← update_last_updated_buffer(pkt.index,active_buffer);
/* for a regular packet */
10 if pkt.type = SKETCH_UPDATE then
/* is this the first time this buffer has been touched since we took a snapshot? */
11 | if active_buffer ≠ last_updated_buffer_for_index then
/* if so, copy data from the inactive buffer before updating */
12 | | if active_buffer = 0 then
13 | | | pkt.result ← copy_update_and_read_buffer_0(pkt.index, pkt.update);
14 | | | else
15 | | | pkt.result ← copy_update_and_read_buffer_1(pkt.index, pkt.update);
/* if not, some other packet has touched this buffer since we took a snapshot, so just
   do update */
16 | | else
17 | | | if active_buffer = 0 then
18 | | | | pkt.result ← update_and_read_buffer_0(pkt.index,pkt.update);
19 | | | | else
20 | | | | pkt.result ← update_and_read_buffer_1(pkt.index, pkt.update);
/* for a snapshot read packet */
21 else if pkt.type = SNAPSHOT_READ then
22 | pkt.update = 0; /* is this the first time this buffer has been touched since we took a
   snapshot? */
23 | if active_buffer ≠ last_updated_buffer_for_index then
/* if so, copy data from the inactive buffer before updating */
24 | | if active_buffer = 0 then
25 | | | pkt.result ← copy_update_and_read_buffer_0(pkt.index, pkt.update);
26 | | | else
27 | | | pkt.result ← copy_update_and_read_buffer_1(pkt.index, pkt.update);
/* if not, some other packet has touched this buffer since we took a snapshot, so just
   do read */
28 | | else
29 | | | if active_buffer = 0 then
30 | | | | pkt.result ← update_and_read_buffer_1(pkt.index, pkt.update);
31 | | | | else
32 | | | | pkt.result ← update_and_read_buffer_0(pkt.index, pkt.update);
```

Appendix C

TLA+ Specification of RedPlane Protocol

We write a TLA+ specification of RedPlane protocol to model-check its correctness.

EXTENDS *Integers, Sequences, TLC, FiniteSets*
 CONSTANTS *NULL, SWITCHES, LEASE_PERIOD, TOTAL_PKTS*

VARIABLES *query, request_queue, SwitchPacketQueue, RemainingLeasePeriod,*
owner, up, active, AliveNum, global_seqnum, pc

$Exists(val) \triangleq val \neq NULL$
 $RequestingSwitches \triangleq \{sw \in SWITCHES : Exists(query[sw]) \wedge query[sw].type = \text{"request"}\}$

VARIABLES *switch, q, seqnum, round, upSwitches, sent_pkts*

$vars \triangleq \langle query, request_queue, SwitchPacketQueue, RemainingLeasePeriod,$
owner, up, active, AliveNum, global_seqnum, pc, switch, q, seqnum,
round, upSwitches, sent_pkts \rangle

$ProcSet \triangleq \{\text{"StateStore"}\} \cup (SWITCHES) \cup \{\text{"LeaseTimer"}\} \cup \{\text{"pktgen"}\}$

$Init \triangleq$ Global variables
 $\wedge query = [sw \in SWITCHES \mapsto NULL]$
 $\wedge request_queue = \langle \rangle$
 $\wedge SwitchPacketQueue = [sw \in SWITCHES \mapsto 0]$
 $\wedge RemainingLeasePeriod = [sw \in SWITCHES \mapsto 0]$
 $\wedge owner = NULL$
 $\wedge up = [sw \in SWITCHES \mapsto TRUE]$
 $\wedge active = [sw \in SWITCHES \mapsto FALSE]$
 $\wedge AliveNum = Cardinality(SWITCHES)$
 $\wedge global_seqnum = 0$

$\wedge switch = NULL$
 $\wedge q = NULL$

$\wedge seqnum = [self \in SWITCHES \mapsto 0]$
 $\wedge round = [self \in SWITCHES \mapsto 0]$

$\wedge upSwitches = \{\}$
 $\wedge sent_pkts = 0$
 $\wedge pc = [self \in ProcSet \mapsto \text{CASE } self = \text{"StateStore"} \rightarrow \text{"START_STORE"}$
 $\square self \in SWITCHES \rightarrow \text{"START_SWITCH"}$
 $\square self = \text{"LeaseTimer"} \rightarrow \text{"START_TIMER"}$
 $\square self = \text{"pktgen"} \rightarrow \text{"START_PKTGEN"}]$

$START_STORE \triangleq \wedge pc[\text{"StateStore"}] = \text{"START_STORE"}$

$\wedge pc' = [pc \text{ EXCEPT } ![\text{"StateStore"}] = \text{"STORE_PROCESSING"}]$
 $\wedge \text{UNCHANGED } \langle query, request_queue, SwitchPacketQueue,$
RemainingLeasePeriod, owner, up, active,
AliveNum, global_seqnum, switch, q, seqnum,
round, upSwitches, sent_pkts\rangle

$STORE_PROCESSING \triangleq \wedge pc[\text{"StateStore"}] = \text{"STORE_PROCESSING"}$
 $\wedge \text{IF } request_queue \neq \langle \rangle$
 $\text{THEN } \wedge switch' = \text{Head}(request_queue)$
 $\wedge request_queue' = \text{Tail}(request_queue)$
 $\wedge q' = query[switch']$
 $\wedge \text{IF } q'.lease_request = \text{"new"}$
 $\text{THEN } \wedge \text{IF } owner \neq \text{NULL}$
 $\text{THEN } \wedge pc' = [pc \text{ EXCEPT } ![\text{"StateStore"}] = \text{"BUFFERING"}]$
 $\text{ELSE } \wedge pc' = [pc \text{ EXCEPT } ![\text{"StateStore"}] = \text{"TRANSFER_LEASE"}]$
 $\text{ELSE } \wedge \text{IF } q'.lease_request = \text{"renew"}$
 $\text{THEN } \wedge pc' = [pc \text{ EXCEPT } ![\text{"StateStore"}] = \text{"RENEW_LEASE"}]$
 $\text{ELSE } \wedge pc' = [pc \text{ EXCEPT } ![\text{"StateStore"}] = \text{"START_STORE"}]$
 $\text{ELSE } \wedge pc' = [pc \text{ EXCEPT } ![\text{"StateStore"}] = \text{"START_STORE"}]$
 $\wedge \text{UNCHANGED } \langle request_queue, switch, q\rangle$
 $\wedge \text{UNCHANGED } \langle query, SwitchPacketQueue,$
RemainingLeasePeriod, owner, up, active,
AliveNum, global_seqnum, seqnum, round,
upSwitches, sent_pkts\rangle

$TRANSFER_LEASE \triangleq \wedge pc[\text{"StateStore"}] = \text{"TRANSFER_LEASE"}$
 $\wedge query' = [query \text{ EXCEPT } ![switch] = [type \mapsto \text{"response"}] @@ ([last_seqnum \mapsto global_seqnum])]$
 $\wedge RemainingLeasePeriod' = [RemainingLeasePeriod \text{ EXCEPT } ![switch] = LEASE_PERIOD]$
 $\wedge owner' = switch$
 $\wedge pc' = [pc \text{ EXCEPT } ![\text{"StateStore"}] = \text{"START_STORE"}]$
 $\wedge \text{UNCHANGED } \langle request_queue, SwitchPacketQueue, up, active,$
AliveNum, global_seqnum, switch, q, seqnum,
round, upSwitches, sent_pkts\rangle

$BUFFERING \triangleq \wedge pc[\text{"StateStore"}] = \text{"BUFFERING"}$
 $\wedge request_queue' = \text{Append}(request_queue, switch)$
 $\wedge pc' = [pc \text{ EXCEPT } ![\text{"StateStore"}] = \text{"STORE_PROCESSING"}]$
 $\wedge \text{UNCHANGED } \langle query, SwitchPacketQueue, RemainingLeasePeriod,$
owner, up, active, AliveNum, global_seqnum,
switch, q, seqnum, round, upSwitches, sent_pkts\rangle

$RENEW_LEASE \triangleq \wedge pc[\text{"StateStore"}] = \text{"RENEW_LEASE"}$
 $\wedge global_seqnum' = q.write_seq$
 $\wedge query' = [query \text{ EXCEPT } ![switch] = [type \mapsto \text{"response"}] @@ ([last_seqnum \mapsto global_seqnum'])]$
 $\wedge RemainingLeasePeriod' = [RemainingLeasePeriod \text{ EXCEPT } ![switch] = LEASE_PERIOD]$

$\wedge \text{owner}' = \text{switch}$
 $\wedge \text{pc}' = [\text{pc} \text{ EXCEPT } ![\text{"StateStore"}] = \text{"START_STORE"}]$
 $\wedge \text{UNCHANGED} \langle \text{request_queue}, \text{SwitchPacketQueue}, \text{up}, \text{active},$
 $\text{AliveNum}, \text{switch}, \text{q}, \text{seqnum}, \text{round}, \text{upSwitches},$
 $\text{sent_pkts} \rangle$

$\text{statestore} \triangleq \text{START_STORE} \vee \text{STORE_PROCESSING} \vee \text{TRANSFER_LEASE}$
 $\vee \text{BUFFERING} \vee \text{RENEW_LEASE}$

$\text{START_SWITCH}(\text{self}) \triangleq \wedge \text{pc}[\text{self}] = \text{"START_SWITCH"}$
 $\wedge \vee \wedge (\text{up}[\text{self}] \wedge \text{SwitchPacketQueue}[\text{self}] > 0)$
 $\wedge \text{active}' = [\text{active} \text{ EXCEPT } ![\text{self}] = \text{TRUE}]$
 $\wedge \text{IF } \text{RemainingLeasePeriod}[\text{self}] = 0$
 $\text{THEN } \wedge \text{pc}' = [\text{pc} \text{ EXCEPT } ![\text{self}] = \text{"NO_LEASE"}]$
 $\text{ELSE } \wedge \text{pc}' = [\text{pc} \text{ EXCEPT } ![\text{self}] = \text{"HAS_LEASE"}]$
 $\vee \wedge \text{pc}' = [\text{pc} \text{ EXCEPT } ![\text{self}] = \text{"SW_FAILURE"}]$
 $\wedge \text{UNCHANGED } \text{active}$
 $\wedge \text{UNCHANGED} \langle \text{query}, \text{request_queue}, \text{SwitchPacketQueue},$
 $\text{RemainingLeasePeriod}, \text{owner}, \text{up},$
 $\text{AliveNum}, \text{global_seqnum}, \text{switch}, \text{q},$
 $\text{seqnum}, \text{round}, \text{upSwitches}, \text{sent_pkts} \rangle$

$\text{NO_LEASE}(\text{self}) \triangleq \wedge \text{pc}[\text{self}] = \text{"NO_LEASE"}$
 $\wedge \text{query}' = [\text{query} \text{ EXCEPT } ![\text{self}] = [\text{type} \mapsto \text{"request"}] \text{@@} ([\text{lease_request} \mapsto \text{"new"}])]$
 $\wedge \text{request_queue}' = \text{Append}(\text{request_queue}, \text{self})$
 $\wedge \text{pc}' = [\text{pc} \text{ EXCEPT } ![\text{self}] = \text{"WAIT_LEASE_RESPONSE"}]$
 $\wedge \text{UNCHANGED} \langle \text{SwitchPacketQueue}, \text{RemainingLeasePeriod},$
 $\text{owner}, \text{up}, \text{active}, \text{AliveNum}, \text{global_seqnum},$
 $\text{switch}, \text{q}, \text{seqnum}, \text{round}, \text{upSwitches},$
 $\text{sent_pkts} \rangle$

$\text{WAIT_LEASE_RESPONSE}(\text{self}) \triangleq \wedge \text{pc}[\text{self}] = \text{"WAIT_LEASE_RESPONSE"}$
 $\wedge \text{query}[\text{self}].\text{type} = \text{"response"}$
 $\wedge \text{seqnum}' = [\text{seqnum} \text{ EXCEPT } ![\text{self}] = \text{query}[\text{self}].\text{last_seqnum}]$
 $\wedge \text{query}' = [\text{query} \text{ EXCEPT } ![\text{self}] = \text{NULL}]$
 $\wedge \text{pc}' = [\text{pc} \text{ EXCEPT } ![\text{self}] = \text{"HAS_LEASE"}]$
 $\wedge \text{UNCHANGED} \langle \text{request_queue}, \text{SwitchPacketQueue},$
 $\text{RemainingLeasePeriod}, \text{owner}, \text{up},$
 $\text{active}, \text{AliveNum}, \text{global_seqnum},$
 $\text{switch}, \text{q}, \text{round}, \text{upSwitches},$
 $\text{sent_pkts} \rangle$

$\text{HAS_LEASE}(\text{self}) \triangleq \wedge \text{pc}[\text{self}] = \text{"HAS_LEASE"}$
 $\wedge \text{seqnum}' = [\text{seqnum} \text{ EXCEPT } ![\text{self}] = \text{seqnum}[\text{self}] + 1]$
 $\wedge \text{query}' = [\text{query} \text{ EXCEPT } ![\text{self}] = [\text{type} \mapsto \text{"request"}] \text{@@} ([\text{lease_request} \mapsto \text{"renew"}, \text{write_seq} \mapsto \text{seqnum}'[\text{self}])]$

$\wedge request_queue' = Append(request_queue, self)$
 $\wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"WAIT_WRITE_RESPONSE"}]$
 $\wedge \text{UNCHANGED } \langle SwitchPacketQueue, RemainingLeasePeriod,$
 $owner, up, active, AliveNum, global_seqnum,$
 $switch, q, round, upSwitches, sent_pkts \rangle$

$WAIT_WRITE_RESPONSE(self) \triangleq \wedge pc[self] = \text{"WAIT_WRITE_RESPONSE"}$
 $\wedge query[self].type = \text{"response"}$
 $\wedge Assert(seqnum[self] = query[self].last_seqnum,$
 $\text{"assertion failed."})$
 $\wedge query' = [query \text{ EXCEPT } ![self] = NULL]$
 $\wedge active' = [active \text{ EXCEPT } ![self] = FALSE]$
 $\wedge SwitchPacketQueue' = [SwitchPacketQueue \text{ EXCEPT } ![self] = SwitchPacketQueue[self] - 1]$
 $\wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"START_SWITCH"}]$
 $\wedge \text{UNCHANGED } \langle request_queue,$
 $RemainingLeasePeriod, owner, up,$
 $AliveNum, global_seqnum, switch,$
 $q, seqnum, round, upSwitches,$
 $sent_pkts \rangle$

$SW_FAILURE(self) \triangleq \wedge pc[self] = \text{"SW_FAILURE"}$
 $\wedge \text{IF } AliveNum > 1 \wedge up[self] = \text{TRUE}$
 $\text{THEN } \wedge up' = [up \text{ EXCEPT } ![self] = FALSE]$
 $\wedge AliveNum' = AliveNum - 1$
 $\wedge query' = query$
 $\text{ELSE } \wedge \text{IF } up[self] = \text{FALSE}$
 $\text{THEN } \wedge up' = [up \text{ EXCEPT } ![self] = \text{TRUE}]$
 $\wedge query' = [query \text{ EXCEPT } ![self] = NULL]$
 $\wedge AliveNum' = AliveNum + 1$
 $\text{ELSE } \wedge \text{TRUE}$
 $\wedge \text{UNCHANGED } \langle query, up, AliveNum \rangle$
 $\wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"START_SWITCH"}]$
 $\wedge \text{UNCHANGED } \langle request_queue, SwitchPacketQueue,$
 $RemainingLeasePeriod, owner, active,$
 $global_seqnum, switch, q, seqnum, round,$
 $upSwitches, sent_pkts \rangle$

$switch_(self) \triangleq START_SWITCH(self) \vee NO_LEASE(self)$
 $\vee WAIT_LEASE_RESPONSE(self) \vee HAS_LEASE(self)$
 $\vee WAIT_WRITE_RESPONSE(self) \vee SW_FAILURE(self)$

$START_TIMER \triangleq \wedge pc[\text{"LeaseTimer"}] = \text{"START_TIMER"}$
 $\wedge owner \neq NULL$
 $\wedge \text{IF } RemainingLeasePeriod[owner] > 0 \wedge active[owner] = \text{FALSE}$
 $\text{THEN } \wedge RemainingLeasePeriod' = [RemainingLeasePeriod \text{ EXCEPT } ![owner] = RemainingLeasePeriod[owner] - 1]$

$\wedge owner' = owner$
 ELSE \wedge IF $RemainingLeasePeriod[owner] = 0$
 THEN $\wedge owner' = NULL$
 ELSE \wedge TRUE
 $\wedge owner' = owner$
 \wedge UNCHANGED $RemainingLeasePeriod$
 $\wedge pc' = [pc \text{ EXCEPT } !["LeaseTimer"]] = \text{"START_TIMER"}$
 \wedge UNCHANGED $\langle query, request_queue, SwitchPacketQueue, up,$
 $active, AliveNum, global_seqnum, switch, q,$
 $seqnum, round, upSwitches, sent_pkts \rangle$

$expirationTimer \triangleq START_TIMER$

$START_PKTGEN \triangleq \wedge pc["pktgen"] = \text{"START_PKTGEN"}$
 \wedge IF $sent_pkts < TOTAL_PKTS$
 THEN $\wedge AliveNum \geq 1$
 $\wedge upSwitches' = \{sw \in SWITCHES : up[sw]\}$
 $\wedge \exists sw \in upSwitches' :$
 $SwitchPacketQueue' = [SwitchPacketQueue \text{ EXCEPT } ![sw] = SwitchPacketQueue[sw] + 1]$
 $\wedge sent_pkts' = sent_pkts + 1$
 $\wedge pc' = [pc \text{ EXCEPT } !["pktgen"]] = \text{"START_PKTGEN"}$
 ELSE $\wedge pc' = [pc \text{ EXCEPT } !["pktgen"]] = \text{"Done"}$
 \wedge UNCHANGED $\langle SwitchPacketQueue, upSwitches,$
 $sent_pkts \rangle$
 \wedge UNCHANGED $\langle query, request_queue, RemainingLeasePeriod,$
 $owner, up, active, AliveNum, global_seqnum,$
 $switch, q, seqnum, round \rangle$

$packetGen \triangleq START_PKTGEN$

$Next \triangleq statestore \vee expirationTimer \vee packetGen$
 $\vee (\exists self \in SWITCHES : switch_ (self))$

$Spec \triangleq \wedge Init \wedge \square [Next]_{vars}$
 $\wedge WF_{vars}(statestore)$
 $\wedge \forall self \in SWITCHES : WF_{vars}(switch_ (self))$
 $\wedge WF_{vars}(expirationTimer)$
 $\wedge WF_{vars}(packetGen)$

$AtLeastOneAliveSwitch \triangleq$
 $\wedge AliveNum \geq 1$
 $\wedge \exists sw \in SWITCHES : up[sw] = \text{TRUE}$

$SingleOwnerInvariant \triangleq$

Bibliography

- [1] Intel ixp2800 network processor. http://www.ic72.com/pdf_file/i/587106.pdf, 2001.
- [2] 2009-M57-Patents packet trace. <http://downloads.digitalcorpora.org/corpora/scenarios/2009-m57-patents/net/>, 2009.
- [3] Data Set for IMC 2010 Data Center Measurement. http://pages.cs.wisc.edu/~tbenson/IMC10_Data.html, 2010.
- [4] pktgen-dpdk: Traffic generator powered by DPDK. <https://git.dpdk.org/apps/pktgen-dpdk/>, 2011.
- [5] RDMA Aware Networks Programming User Manual. https://www.mellanox.com/related-docs/prod_software/RDMA_Aware_Programming_user_manual.pdf, 2011.
- [6] 802.1qbb – priority-based flow control. <https://1.ieee802.org/dcb/802-1qbb/>, 2011.
- [7] A signaling storm is gathering - is your packet core ready? <https://www.nokia.com/blog/a-signaling-storm-is-gathering-is-your-packet-core-ready/>, 2012.
- [8] Intel ethernet switch fm6000 series. <https://www.intel.com/content/dam/www/public/us/en/documents/product-briefs/ethernet-switch-fm6000-series-brief.pdf>, 2014.
- [9] Open network operating system (onos) sdn controller for sdn/nfv solutions. <https://opennetworking.org/onos/>, 2014.
- [10] Intel Xeon Processor E5-2640 v3. <https://ark.intel.com/content/www/us/en/ark/products/83359/intel-xeon-processor-e5-2640-v3-20m-cache-2-60-ghz.html>, 2015.
- [11] vEPC Acceleration Using Agilio SmartNICs. https://www.netronome.com/media/documents/SB_vEPC.pdf, 2017.
- [12] Barefoot Networks Unveils Tofino 2, the Next Generation of the World’s First Fully P4-Programmable Network Switch ASICs. <https://www.barefootnetworks.com/press-releases/barefoot-networks-unveils-tofino-2-the-next-generation-of-the-worlds-first-fully-p4-programmable-network-switch-asics/>, 2018.
- [13] Advanced network telemetry. <https://www.barefootnetworks.com/use-cases/advanced-network-telemetry/>, 2018.
- [14] iperf3. <http://software.es.net/iperf/>, 2018.

- [15] BCM88690–10 Tb/s StrataDNX Jericho2 Ethernet Switch Series. <https://www.broadcom.com/products/ethernet-connectivity/switching/stratadnx/bcm88690>, 2018.
- [16] Agilio CX SmartNICs - Netronome. <https://www.netronome.com/products/agilio-cx/>, 2018.
- [17] Offloading vnfs to programmable switches using p4. <https://wiki.onosproject.org/download/attachments/12420314/p4-vnf-offloading-ons2018.pdf>, 2018.
- [18] Perfctest package. <https://github.com/linux-rdma/perfctest>, 2018.
- [19] P4-SDNet User Guide. https://www.xilinx.com/support/documentation/sw_manuals/xilinx2017_4/ug1252-p4-sdnet.pdf, 2018.
- [20] Vpp - fd.io. <https://wiki.fd.io/view/VPP>, 2018.
- [21] Cavium xpliant ethernet switches. <https://www.cavium.com/xpliant-ethernet-switch-product-family.html>, 2018.
- [22] Amazon ec2 f1 instances. <https://aws.amazon.com/ec2/instance-types/f1/>, 2019.
- [23] Azure VPN Gateway. <https://docs.microsoft.com/en-us/azure/vpn-gateway/vpn-gateway-about-vpngateways>, 2019.
- [24] In-network DDoS Detection. <https://www.barefootnetworks.com/use-cases/in-nw-DDoS-detection/>, 2019.
- [25] Cisco Visual Networking Index. <https://www.cisco.com/c/en/us/solutions/collateral/service-provider/visual-networking-index-vni/white-paper-c11-738429.html>, 2019.
- [26] Fastclick. <https://github.com/tbarbette/fastclick>, 2019.
- [27] Mellanox innova-2 flex open programmable smartnic. http://www.mellanox.com/related-docs/prod_adapter_cards/PB_Innova-2_Flex.pdf, 2019.
- [28] Compare Kemp LoadMaster, F5 Big-IP & Citrix Netscaler. <https://kemptechnologies.com/compare-kemp-f5-big-ip-citrix-netscaler-hardware-load-balancers/>, 2019.
- [29] Liquidio ii 10/25gbe adapter family. <https://www.marvell.com/ethernet-adapters-and-controllers/liquidio-smart-nics/liquidio-ii-smart-nics/index.jsp>, 2019.
- [30] Netronome agilio smartnics. <https://www.netronome.com/products/smartnic/overview/>, 2019.
- [31] P4₁₆ Language Specification. <https://p4.org/p4-spec/docs/P4-16-v1.2.0.html>, 2019.
- [32] EdgeCore Wedge 100BF-32X. <https://www.edge-core.com/productsInfo.php?cls=1&cls2=5&cls3=181&id=335>, 2019.
- [33] Arista 7170 Series. <https://www.arista.com/en/products/7170-series>, 2020.
- [34] NVIDIA Bluefield Data Processing Units. <https://www.nvidia.com/en-us/networking/products/data-processing-unit/>, 2020.

- [35] NVIDIA Mellanox ConnectX-6 Dx. <https://www.nvidia.com/en-us/networking/ethernet/connectx-6-dx/>, 2020.
- [36] The Software Switch Pipeline. https://doc.dpdk.org/guides/prog_guide/packet_framework.html#the-software-switch-swx-pipeline, 2020.
- [37] The Evolved Packet Core. <https://www.3gpp.org/technologies/keywords-acronyms/100-the-evolved-packet-core>, 2020.
- [38] Floodlight SDN OpenFlow Controller. <https://github.com/floodlight/floodlight>, 2020.
- [39] Gurobi - C++ API Overview. https://www.gurobi.com/documentation/9.1/refman/cpp_api_overview.html, 2020.
- [40] iperf(1) - linux man page. <https://linux.die.net/man/1/iperf>, 2020.
- [41] Intel FPGA Programmable Acceleration Card N3000. https://www.intel.com/content/www/us/en/programmable/products/boards_and_kits/dev-kits/altera/intel-fpga-pac-n3000/overview.html, 2020.
- [42] Netcope P4. <https://www.intel.com/content/www/us/en/programmable/solutions/partners/partner-profile/netcope-technologies--a-s-/ip/netcope-p4.html>, 2020.
- [43] NPL Specifications. <https://nplang.org/npl/specifications/>, 2020.
- [44] P4 DPDK backend. <https://github.com/p4lang/p4c/tree/master/backends/dpdk>, 2020.
- [45] Barefoot P4 Studio. <https://www.barefootnetworks.com/products/brief-p4-studio/>, 2020.
- [46] p4c: a reference P4 compiler. <https://github.com/p4lang/p4c>, 2020.
- [47] Pensando DSC-25 Distributed Services Card. <https://pensando.io/wp-content/uploads/2020/03/Pensando-DSC-25-Product-Brief.pdf>, 2020.
- [48] Redis. <https://redis.io/>, 2020.
- [49] Recommended network configuration examples for roce deployment. <https://community.mellanox.com/s/article/recommended-network-configuration-examples-for-roce-deployment>, 2020.
- [50] Cisco Silicon One Q200 and Q200L Processors Data Sheet. <https://www.cisco.com/c/en/us/solutions/collateral/silicon-one/datasheet-c78-744312.html>, 2020.
- [51] Apache Thrift. <https://thrift.apache.org/>, 2020.
- [52] The TLA+ Home Page. <https://lamport.azurewebsites.net/tla/tla.html>, 2020.
- [53] Intel Tofino. <https://www.intel.com/content/www/us/en/products/network-io/programmable-ethernet-switch/tofino-series/tofino.html>, 2020.
- [54] Trident4 / BCM56880 Series. <https://www.broadcom.com/products/ethernet-connectivity/switching/strataxgs/bcm56880-series>, 2020.

- [55] Alveo U280 Data Center Accelerator Card. <https://www.xilinx.com/products/boards-and-kits/alveo/u280.html>, 2020.
- [56] Alveo U50 Data Center Accelerator Card. <https://www.xilinx.com/products/boards-and-kits/alveo/u50.html>, 2020.
- [57] The CAIDA UCSD Anonymized Internet Traces. https://www.caida.org/data/passive/passive_dataset.xml, 2021.
- [58] Netfpga. <https://netfpga.org/index.html>, 2021.
- [59] Redplane public repository. <https://github.com/daehyeok-kim/redplane-public>, 2021.
- [60] 3GPP. 3GPP TS 23.007: Restoration procedures, 2012.
- [61] Albert Gran Alcoz, Alexander Dietmüller, and Laurent Vanbever. Sp-pifo: approximating push-in first-out behaviors using strict-priority queues. In *USENIX NSDI*, 2020.
- [62] D. S. Alexander, W. A. Arbaugh, M. W. Hicks, P. Kakkar, A. D. Keromytis, J. T. Moore, C. A. Gunter, S. M. Nettles, and J. M. Smith. The switchware active network architecture. *Network. Mag. of Global Internetwkg.*, 12(3):29–36, May 1998.
- [63] Thomas Anderson and Michael Dahlin. *Operating Systems: Principles and Practice*, volume 1. Recursive books.
- [64] Alexey Andreyev. Introducing data center fabric, the next-generation facebook data center network. <https://engineering.fb.com/production-engineering/introducing-data-center-fabric-the-next-generation-facebook-data-center-network/>, 2014.
- [65] Mina Tahmasbi Arashloo, Pavel Shirshov, Rohan Gandhi, Guohan Lu, Lihua Yuan, and Jennifer Rexford. A scalable vpn gateway for multi-tenant cloud services. *ACM SIGCOMM Computer Communication Review*, 48(1):49–55, April 2018.
- [66] Barefoot Networks. Tofino Switch Architecture Specification (accessible under NDA), 2017.
- [67] Theophilus Benson, Aditya Akella, and David A. Maltz. Network traffic characteristics of data centers in the wild. In *ACM IMC*, 2010.
- [68] Philip A. Bernstein, Vassos Hadzilacos, and Nathan Goodman. *Concurrency Control and Recovery in Database Systems*. Addison-Wesley.
- [69] Samrat Bhattacharjee, Kenneth L Calvert, and Ellen W Zegura. An architecture for active networking. In *International Conference on High Performance Networking*, pages 265–279, 1997.
- [70] Burton H. Bloom. Space/time trade-offs in hash coding with allowable errors. *Commun. ACM*, 13(7):422–426, July 1970.
- [71] Pat Bosshart, Glen Gibb, Hun-Seok Kim, George Varghese, Nick McKeown, Martin Izzard, Fernando Mujica, and Mark Horowitz. Forwarding metamorphosis: Fast programmable match-action processing in hardware for sdn. In *ACM SIGCOMM*, 2013.
- [72] Pat Bosshart, Dan Daly, Glen Gibb, Martin Izzard, Nick McKeown, Jennifer Rexford, Cole Schlesinger, Dan Talayco, Amin Vahdat, George Varghese, and David Walker. P4: Program-

- ming protocol-independent packet processors. *ACM SIGCOMM Computer Communication Review*, 44(3):87–95, 2014.
- [73] K. L. Calvert, S. Bhattacharjee, E. Zegura, and J. Sterbenz. Directions in active networks. *IEEE Communications Magazine*, 36(10):72–78, Oct 1998.
- [74] Brian Carpenter and Scott Brim. Middleboxes: Taxonomy and issues. Technical report, RFC 3234, February, 2002.
- [75] Martin Casado, Tal Garfinkel, Aditya Akella, Michael J Freedman, Dan Boneh, Nick McKeown, and Scott Shenker. Sane: A protection architecture for enterprise networks. In *USENIX Security*, 2006.
- [76] Martin Casado, Michael J Freedman, Justin Pettit, Jianying Luo, Nick McKeown, and Scott Shenker. Ethane: Taking control of the enterprise. *ACM SIGCOMM Computer Communication Review*, 37(4):1–12, 2007.
- [77] Adrian M. Caulfield, Eric S. Chung, Andrew Putnam, Hari Angepat, Jeremy Fowers, Michael Haselman, Stephen Heil, Matt Humphrey, Puneet Kaur, Joo-Young Kim, Daniel Lo, Todd Massengill, Kalin Ovtcharov, Michael Papamichael, Lisa Woods, Sitaram Lanka, Derek Chiou, and Doug Burger. A cloud-scale acceleration architecture. In *IEEE/ACM MICRO*, 2016.
- [78] Xiaoqi Chen, Shir Landau-Feibish, Mark Braverman, and Jennifer Rexford. Beaucoup: Answering many network traffic queries, one memory update at a time. In *ACM SIGCOMM*, 2020.
- [79] Yanzhe Chen, Xingda Wei, Jiabin Shi, Rong Chen, and Haibo Chen. Fast and general distributed transactions using rdma and htm. In *EuroSys*, 2016.
- [80] Sharad Chole, Andy Fingerhut, Sha Ma, Anirudh Sivaraman, Shay Vargaftik, Alon Berger, Gal Mendelson, Mohammad Alizadeh, Shang-Tse Chuang, Isaac Keslassy, et al. drmt: Disaggregated programmable switching. In *ACM SIGCOMM*, 2017.
- [81] Cisco. Cisco Global Cloud Index: Forecast and Methodology, 2016–2021 White Paper.
- [82] David Clark. The design philosophy of the darpa internet protocols. In *Symposium proceedings on Communications architectures and protocols*, 1988.
- [83] David Clayton, Christopher Patton, and Thomas Shrimpton. Probabilistic data structures in adversarial environments. In *ACM CCS*, 2019.
- [84] Graham Cormode and Marios Hadjieleftheriou. Finding frequent items in data streams. *Proceedings of the VLDB Endowment*, 1(2), 2008.
- [85] Michael Dalton, David Schultz, Jacob Adriaens, Ahsan Arefin, Anshuman Gupta, Brian Fahs, Dima Rubinstein, Enrique Cauich Zermeno, Erik Rubow, James Alexander Docauer, Jesse Alpert, Jing Ai, Jon Olson, Kevin DeCabooter, Marc de Kruijf, Nan Hua, Nathan Lewis, Nikhil Kasinadhuni, Riccardo Crepaldi, Srinivas Krishnan, Subbaiah Venkata, Yossi Richter, Uday Naik, and Amin Vahdat. Andromeda: Performance, isolation, and velocity at scale in cloud network virtualization. In *USENIX NSDI*, 2018.

- [86] H. T. Dang, P. Bressana, H. Wang, K. S. Lee, N. Zilberman, H. Weatherspoon, M. Canini, F. Pedone, and R. Soulé. P4xos: Consensus as a network service. *IEEE/ACM Transactions on Networking*, 2020.
- [87] Mihai Dobrescu, Norbert Egi, Katerina Argyraki, Byung-Gon Chun, Kevin Fall, Gianluca Iannaccone, Allan Knies, Maziar Manesh, and Sylvia Ratnasamy. Routebricks: Exploiting parallelism to scale software routers. In *ACM SOSP*, 2009.
- [88] Aleksandar Dragojević, Dushyanth Narayanan, Miguel Castro, and Orion Hodson. Farm: Fast remote memory. In *USENIX NSDI*, 2014.
- [89] Aleksandar Dragojević, Dushyanth Narayanan, Edmund B. Nightingale, Matthew Renzelmann, Alex Shamis, Anirudh Badam, and Miguel Castro. No compromises: Distributed transactions with consistency, availability, and performance. In *ACM SOSP*, 2015.
- [90] Daniel E. Eisenbud, Cheng Yi, Carlo Contavalli, Cody Smith, Roman Kononov, Eric Mann-Hielscher, Ardas Cilingiroglu, Bin Cheyney, Wentao Shang, and Jinnah Dylan Hosein. Maglev: A fast and reliable software network load balancer. In *USENIX NSDI*, 2015.
- [91] Bin Fan, Hyeontaek Lim, David G. Andersen, and Michael Kaminsky. Small cache, big effect: Provable load balancing for randomly partitioned cluster services. In *ACM SOCC*, 2011.
- [92] Nick Feamster, Hari Balakrishnan, Jennifer Rexford, Aman Shaikh, and Jacobus Van Der Merwe. The case for separating routing from routers. In *ACM SIGCOMM workshop on Future directions in network architecture*, 2004.
- [93] Nick Feamster, Jennifer Rexford, and Ellen Zegura. The road to sdn: an intellectual history of programmable networks. *ACM SIGCOMM Computer Communication Review*, 44(2):87–98, 2014.
- [94] Daniel Firestone, Andrew Putnam, Sambhrama Mundkur, Derek Chiou, Alireza Dabagh, Mike Andrewartha, Hari Angepat, Vivek Bhanu, Adrian Caulfield, Eric Chung, Harish Kumar Chandrappa, Somesh Chaturmohta, Matt Humphrey, Jack Lavier, Norman Lam, Fengfen Liu, Kalin Ovtcharov, Jitu Padhye, Gautham Popuri, Shachar Raindel, Tejas Sapre, Mark Shaw, Gabriel Silva, Madhan Sivakumar, Nisheeth Srivastava, Anshuman Verma, Qasim Zuhair, Deepak Bansal, Doug Burger, Kushagra Vaid, David A. Maltz, and Albert Greenberg. Azure accelerated networking: Smartnics in the public cloud. In *USENIX NSDI*, 2018.
- [95] Rohan Gandhi, Hongqiang Harry Liu, Y. Charlie Hu, Guohan Lu, Jitendra Padhye, Lihua Yuan, and Ming Zhang. Duet: Cloud scale load balancing with hardware and software. In *ACM SIGCOMM*, 2014.
- [96] Jiaqi Gao, Ennan Zhai, Hongqiang Harry Liu, Rui Miao, Yu Zhou, Bingchuan Tian, Chen Sun, Dennis Cai, Ming Zhang, and Minlan Yu. Lyra: A cross-platform language and compiler for data plane programming on heterogeneous asics. In *ACM SIGCOMM*, 2020.

- [97] Aaron Gember-Jacobson, Raajay Viswanathan, Chaithan Prakash, Robert Grandl, Junaid Khalid, Sourav Das, and Aditya Akella. Opennf: Enabling innovation in network function control. In *ACM SIGCOMM*, 2014.
- [98] Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung. The google file system. In *ACM SOSP*, 2003.
- [99] Phillipa Gill, Navendu Jain, and Nachiappan Nagappan. Understanding network failures in data centers: measurement, analysis, and implications. In *ACM SIGCOMM*, 2011.
- [100] Cary Gray and David Cheriton. Leases: An efficient fault-tolerant mechanism for distributed file cache consistency. *ACM SIGOPS Operating Systems Review*, 23(5), 1989.
- [101] Albert Greenberg, Gisli Hjalmytsson, David A Maltz, Andy Myers, Jennifer Rexford, Geoffrey Xie, Hong Yan, Jibin Zhan, and Hui Zhang. A clean slate 4d approach to network control and management. *ACM SIGCOMM Computer Communication Review*, 35(5):41–54, 2005.
- [102] Albert Greenberg, James R Hamilton, Navendu Jain, Srikanth Kandula, Changhoon Kim, Parantap Lahiri, David A Maltz, Parveen Patel, and Sudipta Sengupta. VL2: a scalable and flexible data center network. In *ACM SIGCOMM*, 2009.
- [103] Natasha Gude, Teemu Koponen, Justin Pettit, Ben Pfaff, Martín Casado, Nick McKeown, and Scott Shenker. Nox: towards an operating system for networks. *ACM SIGCOMM Computer Communication Review*, 38(3):105–110, 2008.
- [104] Chuanxiong Guo, Haitao Wu, Zhong Deng, Gaurav Soni, Jianxi Ye, Jitu Padhye, and Marina Lipshteyn. Rdma over commodity ethernet at scale. In *ACM SIGCOMM*, 2016.
- [105] Arpit Gupta, Rob Harrison, Marco Canini, Nick Feamster, Jennifer Rexford, and Walter Willinger. Sonata: Query-driven streaming network telemetry. In *ACM SIGCOMM*, 2018.
- [106] Sangjin Han, Keon Jang, KyoungSoo Park, and Sue Moon. Packetshader: a gpu-accelerated software router. *ACM SIGCOMM Computer Communication Review*, 40(4):195–206, 2010.
- [107] David Hancock and Jacobus Van der Merwe. Hyper4: Using p4 to virtualize the programmable data plane. In *ACM CoNEXT*, 2016.
- [108] Mark Handley, Costin Raiciu, Alexandru Agache, Andrei Voinescu, Andrew W Moore, Gianni Antichi, and Marcin Wójcik. Re-architecting datacenter networks and stacks for low latency and high performance. In *ACM SIGCOMM*, 2017.
- [109] Frederik Hauser, Marco Häberle, Daniel Merling, Steffen Lindner, Vladimir Gurevich, Florian Zeiger, Reinhard Frank, and Michael Menth. A survey on data plane programming with p4: Fundamentals, advances, and applied research. *arXiv preprint arXiv:2101.10632*, 2021.
- [110] J. Heinanen and R. Guerin. A two rate three color marker. Technical report, RFC 2698, February, 1999.
- [111] Maurice P Herlihy and Jeannette M Wing. Linearizability: A correctness condition for concurrent objects. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 12(3), 1990.

- [112] Stephen Ibanez, Gordon Brebner, Nick McKeown, and Noa Zilberman. The p4- > netfpga workflow for line-rate packet processing. In *ACM/SIGDA FPGA*, 2019.
- [113] Infiniband Trade Association. Supplement to InfiniBand architecture specification volume 1 release 1.2.1 annex A16: RDMA over converged ethernet (RoCE), 2010.
- [114] Infiniband Trade Association. Supplement to InfiniBand architecture specification volume 1 release 1.2.1 annex A17: RDMA over converged ethernet (RoCE), 2010.
- [115] Xin Jin, Xiaozhou Li, Haoyu Zhang, Robert Soulé, Jeongkeun Lee, Nate Foster, Changhoon Kim, and Ion Stoica. Netcache: Balancing key-value stores with fast in-network caching. In *ACM SOSP*, 2017.
- [116] Xin Jin, Xiaozhou Li, Haoyu Zhang, Nate Foster, Jeongkeun Lee, Robert Soulé, Changhoon Kim, and Ion Stoica. Netchain: Scale-free sub-rtt coordination. In *USENIX NSDI*, 2018.
- [117] Murad Kablan, Azzam Alsudais, Eric Keller, and Franck Le. Stateless network functions: Breaking the tight coupling of state and processing. In *USENIX NSDI*, 2017.
- [118] Anuj Kalia, Michael Kaminsky, and David G. Andersen. Using rdma efficiently for key-value services. In *ACM SIGCOMM*, 2014.
- [119] David Karger, Eric Lehman, Tom Leighton, Rina Panigrahy, Matthew Levine, and Daniel Lewin. Consistent hashing and random trees: Distributed caching protocols for relieving hot spots on the world wide web. In *STOC*, 1997.
- [120] Georgios P Katsikas, Tom Barbette, Dejan Kostic, Rebecca Steinert, and Gerald Q Maguire Jr. Metron: Nfv service chains at the true speed of the underlying hardware. In *USENIX NSDI*, 2018.
- [121] Naga Katta, Haoyu Zhang, Michael Freedman, and Jennifer Rexford. Ravana: Controller fault-tolerance in software-defined networking. In *ACM SOSR*, 2015.
- [122] Naga Katta, Omid Alipourfard, Jennifer Rexford, and David Walker. Cacheflow: Dependency-aware rule-caching for software-defined networks. In *ACM SOSR*, 2016.
- [123] Naga Katta, Mukesh Hira, Changhoon Kim, Anirudh Sivaraman, and Jennifer Rexford. Hula: Scalable load balancing using programmable data planes. In *ACM SOSR*, 2016.
- [124] Eric Keller, Jennifer Rexford, and Jacobus E van der Merwe. Seamless bgp migration with router grafting. In *USENIX NSDI*, 2010.
- [125] Elie F Kfoury, Jorge Crichigno, and Elias Bou-Harb. An exhaustive survey on p4 programmable data plane switches: Taxonomy, applications, challenges, and future trends. *arXiv preprint arXiv:2102.00643*, 2021.
- [126] Changhoon Kim, Matthew Caesar, Alexandre Gerber, and Jennifer Rexford. Revisiting route caching: The world should be flat. In *PAM*, 2009.
- [127] Daehyeok Kim, Amirsaman Memaripour, Anirudh Badam, Yibo Zhu, Hongqiang Harry Liu, Jitu Padhye, Shachar Raindel, Steven Swanson, Vyas Sekar, and Srinivasan Seshan. Hyperloop: Group-based nic-offloading to accelerate replicated transactions in multi-tenant storage systems. In *ACM SIGCOMM*, 2018.

- [128] Daehyeok Kim, Yibo Zhu, Changhoon Kim, Jeongkeun Lee, and Srinivasan Seshan. Generic external memory for switch data planes. In *ACM HotNets*, 2018.
- [129] Daehyeok Kim, Zaoxing Liu, Yibo Zhu, Changhoon Kim, Jeongkeun Lee, Vyas Sekar, and Srinivasan Seshan. Tea: Enabling state-intensive network functions on programmable switches. In *ACM SIGCOMM*, 2020.
- [130] Daehyeok Kim, Jacob Nelson, Dan R. K. Ports, Vyas Sekar, and Srinivasan Seshan. Red-plane: Enabling fault-tolerant stateful in-switch applications. In *ACM SIGCOMM*, 2021.
- [131] Eddie Kohler, Robert Morris, Benjie Chen, John Jannotti, and M Frans Kaashoek. The click modular router. *ACM Transactions on Computer Systems*, 18(3):263–297, 2000.
- [132] Teemu Koponen, Martin Casado, Natasha Gude, Jeremy Stribling, Leon Poutievski, Min Zhu, Rajiv Ramanathan, Yuichiro Iwata, Hiroaki Inoue, Takayuki Hama, et al. Onix: A distributed control platform for large-scale production networks. In *USENIX OSDI*, 2010.
- [133] TV Lakshman, T Nandagopal, R Ramjee, K Sabnani, and T Woo. The softrouter architecture. In *ACM HotNets*, 2004.
- [134] Petr Lapukhov, Ariff Premji, and Jon Mitchell. Use of bgp for routing in large-scale data centers. *Internet Requests for Comments RFC Editor RFC*, 7938, 2016.
- [135] Collin Lee, Seo Jin Park, Ankita Kejriwal, Satoshi Matsushita, and John K. Ousterhout. Implementing linearizability at large scale and low latency. In *ACM SOSP*, 2015.
- [136] Bojie Li, Kun Tan, Layong Larry Luo, Yanqing Peng, Renqian Luo, Ningyi Xu, Yongqiang Xiong, Peng Cheng, and Enhong Chen. Clicknp: Highly flexible and high performance network processing with reconfigurable hardware. In *ACM SIGCOMM*, 2016.
- [137] Jialin Li, Ellis Michael, Naveen Kr. Sharma, Adriana Szekeres, and Dan R. K. Ports. Just say no to paxos overhead: Replacing consensus with network ordering. In *USENIX OSDI*, 2016.
- [138] Jialin Li, Ellis Michael, and Dan RK Ports. Eris: Coordination-free consistent transactions using in-network concurrency control. In *ACM SOSP*, 2017.
- [139] Jialin Li, Jacob Nelson, Ellis Michael, Xin Jin, and Dan RK Ports. Pegasus: Tolerating skewed workloads in distributed storage with in-network coherence directories. In *USENIX OSDI*, 2020.
- [140] Yuliang Li, Rui Miao, Hongqiang Harry Liu, Yan Zhuang, Fei Feng, Lingbo Tang, Zheng Cao, Ming Zhang, Frank Kelly, Mohammad Alizadeh, et al. Hpcc: High precision congestion control. In *ACM SIGCOMM*. 2019.
- [141] Hyeontaek Lim, Dongsu Han, David G Andersen, and Michael Kaminsky. Mica: A holistic approach to fast in-memory key-value storage. In *USENIX NSDI*, 2014.
- [142] Barbara Liskov and James Cowling. Viewstamped replication revisited. Technical Report MIT-CSAIL-TR-2012-021, MIT Computer Science and Artificial Intelligence Laboratory, Cambridge, MA, USA, July 2012.

- [143] Hongqiang Harry Liu, Yibo Zhu, Jitu Padhye, Jiabin Cao, Sri Tallapragada, Nuno P Lopes, Andrey Rybalchenko, Guohan Lu, and Lihua Yuan. Crystalnet: Faithfully emulating large production networks. In *ACM SOSP*, 2017.
- [144] Ming Liu, Liang Luo, Jacob Nelson, Luis Ceze, Arvind Krishnamurthy, and Kishore Atreya. Incbricks: Toward in-network computation with an in-network cache. In *ACM ASPLOS*, 2017.
- [145] Vincent Liu, Daniel Halperin, Arvind Krishnamurthy, and Thomas Anderson. F10: A fault-tolerant engineered network. In *USENIX NSDI*, 2013.
- [146] Zaoxing Liu, Antonis Manousis, Gregory Vorsanger, Vyas Sekar, and Vladimir Braverman. One sketch to rule them all: Rethinking network flow monitoring with univmon. In *ACM SIGCOMM*, 2016.
- [147] Zaoxing Liu, Zhihao Bai, Zhenming Liu, Xiaozhou Li, Changhoon Kim, Vladimir Braverman, Xin Jin, and Ion Stoica. Distcache: Provable load balancing for large-scale storage systems with distributed caching. In *USENIX FAST*, 2019.
- [148] Zaoxing Liu, Hun Namkung, Georgios Nikolaidis, Jeongkeun Lee, Changhoon Kim, Xin Jin, Vladimir Braverman, Minlan Yu, and Vyas Sekar. Jaqen: A high-performance switch-native approach for detecting and mitigating volumetric ddos attacks with programmable switches. In *USENIX Security*, 2021.
- [149] Kirk Lougheed and Yakov Rekhter. Border gateway protocol (bgp). Technical report, IETF RFC 1105, june, 1989.
- [150] Youyou Lu, Jiwu Shu, Youmin Chen, and Tao Li. Octopus: an rdma-enabled distributed persistent memory file system. In *USENIX ATC*, 2017.
- [151] Nick McKeown, Tom Anderson, Hari Balakrishnan, Guru Parulkar, Larry Peterson, Jennifer Rexford, Scott Shenker, and Jonathan Turner. Openflow: Enabling innovation in campus networks. *ACM SIGCOMM Computer Communication Review*, 38(2):69–74, 2008.
- [152] Justin Meza, Tianyin Xu, Kaushik Veeraraghavan, and Onur Mutlu. A large scale study of data center network reliability. In *ACM IMC*, 2018.
- [153] Rui Miao, Hongyi Zeng, Changhoon Kim, Jeongkeun Lee, and Minlan Yu. Silkroad: Making stateful layer-4 load balancing fast and cheap using switching asics. In *ACM SIGCOMM*, 2017.
- [154] Christopher Mitchell, Yifeng Geng, and Jinyang Li. Using one-sided RDMA reads to build a fast, cpu-efficient key-value store. In *USENIX ATC*, 2013.
- [155] Radhika Mittal, Vinh The Lam, Nandita Dukkipati, Emily Blem, Hassan Wassel, Monia Ghobadi, Amin Vahdat, Yaogong Wang, David Wetherall, and David Zats. Timely: Rtt-based congestion control for the datacenter. In *ACM SIGCOMM*, 2015.
- [156] Ali Mohammadkhan, KK Ramakrishnan, Ashok Sunder Rajan, and Christian Maciocco. Considerations for re-designing the cellular infrastructure exploiting software-based networks. In *IEEE ICNP*, 2016.

- [157] Srinivas Narayana, Anirudh Sivaraman, Vikram Nathan, Prateesh Goyal, Venkat Arun, Mohammad Alizadeh, Vimalkumar Jeyakumar, and Changhoon Kim. Language-directed hardware design for network performance monitoring. In *ACM SIGCOMM*, 2017.
- [158] Diego Ongaro and John Ousterhout. In search of an understandable consensus algorithm. In *USENIX ATC*, 2014.
- [159] Diego Ongaro, Stephen M. Rumble, Ryan Stutsman, John Ousterhout, and Mendel Rosenblum. Fast crash recovery in ramcloud. In *ACM SOSP*, 2011.
- [160] John Ousterhout, Parag Agrawal, David Erickson, Christos Kozyrakis, Jacob Leverich, David Mazières, Subhasish Mitra, Aravind Narayanan, Guru Parulkar, Mendel Rosenblum, et al. The case for ramclouds: scalable high-performance storage entirely in dram. *ACM SIGOPS Operating Systems Review*, 43(4):92–105, 2010.
- [161] Rasmus Pagh and Flemming Friche Rodler. Cuckoo hashing. In *European Symposium on Algorithms*, 2001.
- [162] Tian Pan, Nianbing Yu, Chenhao Jia, Jianwen Pi, Liang Xu, Yisong Qiao, Zhiguo Li, Kun Liu, Jie Lu, Jianyuan Lu, Enge Song, Jiao Zhang, Tao Huang, and Shunmin Zhu. Sailfish: Accelerating cloud-scale multi-tenant multi-service gateways with programmable switches. In *ACM SIGCOMM*, 2021.
- [163] Parveen Patel, Deepak Bansal, Lihua Yuan, Ashwin Murthy, Albert Greenberg, David A. Maltz, Randy Kern, Hemant Kumar, Marios Zikos, Hongyu Wu, Changhoon Kim, and Naveen Karri. Ananta: Cloud scale load balancing. In *ACM SIGCOMM*, 2013.
- [164] Larry Peterson, Tom Anderson, Sachin Katti, Nick McKeown, Guru Parulkar, Jennifer Rexford, Mahadev Satyanarayanan, Oguz Sunay, and Amin Vahdat. Democratizing the network edge. *ACM SIGCOMM Computer Communication Review*, 49(2), 2019.
- [165] S. Pontarelli, P. Reviriego, and M. Mitzenmacher. Emoma: Exact match in one memory access. *IEEE Transactions on Knowledge and Data Engineering*, 30(11):2120–2133, Nov 2018.
- [166] Chaithan Prakash, Jeongkeun Lee, Yoshio Turner, Joon-Myung Kang, Aditya Akella, Sujata Banerjee, Charles Clark, Yadi Ma, Puneet Sharma, and Ying Zhang. Pga: Using graphs to express and automatically reconcile network policies. *ACM SIGCOMM Computer Communication Review*, 45(4):29–42, 2015.
- [167] Zafar Ayyub Qazi, Melvin Walls, Aurojit Panda, Vyas Sekar, Sylvia Ratnasamy, and Scott Shenker. A high performance packet core for next generation cellular networks. In *ACM SIGCOMM*, 2017.
- [168] Kun Qian, Sai Ma, Mao Miao, Jianyuan Lu, Tong Zhang, Peilong Wang, Chenghao Sun, and Fengyuan Ren. Flexgate: High-performance heterogeneous gateway in data centers. In *APNET*, 2019.
- [169] Shriram Rajagopalan, Dan Williams, and Hani Jamjoom. Pico replication: A high availability framework for middleboxes. In *ACM SOCC*, 2013.
- [170] Shriram Rajagopalan, Dan Williams, Hani Jamjoom, and Andrew Warfield. Split/merge: System support for elastic execution in virtual middleboxes. In *USENIX NSDI*, 2013.

- [171] Mark Reitblatt, Nate Foster, Jennifer Rexford, Cole Schlesinger, and David Walker. Abstractions for network update. *ACM SIGCOMM Computer Communication Review*, 42(4): 323–334, 2012.
- [172] Arjun Roy, Hongyi Zeng, Jasmeet Bagga, George Porter, and Alex C. Snoeren. Inside the social network’s (datacenter) network. In *ACM SIGCOMM*, 2015.
- [173] Jerome H Saltzer, David P Reed, and David D Clark. End-to-end arguments in system design. *ACM Transactions on Computer Systems (TOCS)*, 2(4), 1984.
- [174] Amedeo Sapia, Marco Canini, Chen-Yu Ho, Jacob Nelson, Panos Kalnis, Changhoon Kim, Arvind Krishnamurthy, Masoud Moshref, Dan Ports, and Peter Richtarik. Scaling distributed machine learning with in-network aggregation. In *USENIX NSDI*, 2021.
- [175] Karla Saur, Joseph Collard, Nate Foster, Arjun Guha, Laurent Vanbever, and Michael Hicks. Safe and flexible controller upgrades for sdns. In *ACM SOSR*, 2016.
- [176] Fred B. Schneider. Implementing fault-tolerant services using the state machine approach: a tutorial. *ACM Computing Surveys*, 22(4):299–319, 1990.
- [177] Rinku Shah, Vikas Kumar, Mythili Vutukuru, and Purushottam Kulkarni. Turboepc: Leveraging dataplane programmability to accelerate the mobile packet core. In *ACM SOSR*, 2020.
- [178] Yizhou Shan, Shin-Yeh Tsai, and Yiying Zhang. Distributed shared persistent memory. In *ACM SOCC*, 2017.
- [179] Justine Sherry, Peter Xiang Gao, Soumya Basu, Aurojit Panda, Arvind Krishnamurthy, Christian Maciocco, Maziar Manesh, João Martins, Sylvia Ratnasamy, Luigi Rizzo, et al. Rollback-recovery for middleboxes. In *ACM SIGCOMM*, 2015.
- [180] Arjun Singh, Joon Ong, Amit Agarwal, Glen Anderson, Ashby Armistead, Roy Bannon, Seb Boving, Gaurav Desai, Bob Felderman, Paulie Germano, Anand Kanagala, Jeff Provost, Jason Simmons, Eiichi Tanda, Jim Wanderer, Urs Hölzle, Stephen Stuart, and Amin Vahdat. Jupiter rising: A decade of clos topologies and centralized control in google’s datacenter network. In *ACM SIGCOMM*, 2015.
- [181] Jonathan M Smith and Scott M Nettles. Active networking: one view of the past, present, and future. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 34(1):4–18, 2004.
- [182] Jonathan M Smith, Dave J Farber, CA Gunter, SM Nettles, DC Feldmeier, and William D Sincoskie. Switchware: Accelerating network evolution. *Dept. Comput. Inform. Sci., Univ. Pennsylvania, Philadelphia, Tech. Rep. MS-CIS-96-38*, 1996.
- [183] John Sonchack, Adam J Aviv, Eric Keller, and Jonathan M Smith. Turboflow: Information rich flow record generation on commodity switches. In *ACM EuroSys*, 2018.
- [184] Hardik Soni, Myriana Rifai, Praveen Kumar, Ryan Doenges, and Nate Foster. Composing dataplane programs with $\mu p4$. In *ACM SIGCOMM*, 2020.

- [185] Nik Sultana, John Sonchack, Hans Giesen, Isaac Pedisich, Zhaoyang Han, Nishanth Shyamkumar, Shivani Burad, André DeHon, and Boon Thau Loo. Flightplan: Dataplane disaggregation and placement for p4 programs. In *USENIX NSDI*, 2021.
- [186] L. Tang, Q. Huang, and P. P. C. Lee. Spreadsketch: Toward invertible and network-wide detection of superspreaders. In *IEEE INFOCOM*, 2020.
- [187] David L. Tennenhouse and David J. Wetherall. Towards an active network architecture. *ACM SIGCOMM Computer Communication Review*, 26(2):5–17, April 1996. ISSN 0146-4833.
- [188] David G. Thaler and Chinya V. Ravishankar. Using name-based mappings to increase hit rates. *IEEE/ACM Trans. Netw.*, 6(1):1–14, February 1998.
- [189] Muhammad Tirmazi, Ran Ben Basat, Jiaqi Gao, and Minlan Yu. Cheetah: Accelerating database queries with switch pruning. In *ACM SIGMOD*, 2020.
- [190] William Tu, Fabian Ruffy, and Mihai Budiu. Linux network programming with p4. In *Linux Plumb. Conf*, 2018.
- [191] Robbert van Renesse and Fred B. Schneider. Chain replication for supporting high throughput and availability. In *USENIX OSDI*, 2004.
- [192] Marcel Waldvogel, George Varghese, Jon Turner, and Bernhard Plattner. Scalable high speed ip routing lookups. In *ACM SIGCOMM*, 1997.
- [193] Han Wang, Robert Soulé, Huynh Tu Dang, Ki Suh Lee, Vishal Shrivastav, Nate Foster, and Hakim Weatherspoon. P4fpga: A rapid prototyping framework for p4. In *ACM SOSR*, 2017.
- [194] David J Wetherall, John V Guttag, and David L Tennenhouse. Ants: A toolkit for building and dynamically deploying network protocols. In *IEEE Open Architectures and Network Programming*, 1998.
- [195] T. Wolf and J. S. Turner. Design issues for high-performance active routers. *IEEE Journal on Selected Areas in Communications*, 19(3):404–409, March 2001.
- [196] Shinae Woo, Justine Sherry, Sangjin Han, Sue Moon, Sylvia Ratnasamy, and Scott Shenker. Elastic scaling of stateful network functions. In *USENIX NSDI*, 2018.
- [197] Jiarong Xing, Wenqing Wu, and Ang Chen. Ripple: A programmable, decentralized link-flooding defense against adaptive adversaries. In *USENIX Security*, 2021.
- [198] Lily Yang, Ram Dantu, Terry Anderson, and Ram Gopal. Forwarding and control element separation (forces) framework. Technical report, RFC 3746, April, 2004.
- [199] Tong Yang, Jie Jiang, Peng Liu, Qun Huang, Junzhi Gong, Yang Zhou, Rui Miao, Xiaoming Li, and Steve Uhlig. Elastic sketch: Adaptive and fast network-wide measurements. In *ACM SIGCOMM*, 2018.
- [200] Zhuolong Yu, Yiwen Zhang, Vladimir Braverman, Mosharaf Chowdhury, and Xin Jin. Netlock: Fast, centralized lock management using programmable switches. In *ACM SIGCOMM*, 2020.

- [201] Zhuolong Yu, Chuheng Hu, Jingfeng Wu, Xiao Sun, Vladimir Braverman, Mosharaf Chowdhury, Zhenhua Liu, and Xin Jin. Programmable packet scheduling with a single queue. In *ACM SIGCOMM*, 2021.
- [202] Cheng Zhang, Jun Bi, Yu Zhou, Abdul Basit Dogar, and Jianping Wu. Hyperv: A high performance hypervisor for virtualization of the programmable data plane. In *IEEE ICCCN*, 2017.
- [203] Menghao Zhang, Guanyu Li, Shicheng Wang, Chang Liu, Ang Chen, Hongxin Hu, Guofei Gu, Qianqian Li, Mingwei Xu, and Jianping Wu. Poseidon: Mitigating volumetric ddos attacks with programmable switches. In *ISOC NDSS*, 2020.
- [204] Peng Zheng, Theophilus Benson, and Chengchen Hu. P4visor: Lightweight virtualization and composition primitives for building and testing modular programs. In *ACM CoNEXT*, 2018.
- [205] Dong Zhou, Huacheng Yu, Michael Kaminsky, and David Andersen. Fast software cache design for network appliances. In *USENIX ATC*, 2020.
- [206] Hang Zhu, Zhihao Bai, Jialin Li, Ellis Michael, Dan RK Ports, Ion Stoica, and Xin Jin. Harmonia: Near-linear scalability for replicated storage with in-network conflict detection. *Proceedings of the VLDB Endowment*, 13(3), 2019.