

Network Virtualization in a Shared TCP Fast-Path

by

Mehrshad Lotfi Foroushani

Master Thesis

Software Engineering Chair
Faculty of Natural Sciences and Technology I
Department of Computer Science
Saarland University

Supervisor

Antoine Kaufmann

Advisor

Antione Kaufmann

Reviewers

Antione Kaufmann

Anja Feldmann

March 19, 2023



Declaration of Authorship

Eidesstattliche Erklärung

Ich erkläre hiermit an Eides Statt, dass ich die vorliegende Arbeit selbstständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel verwendet habe.

Statement in Lieu of an Oath

I hereby confirm that I have written this thesis on my own and that I have not used any other media or materials than the ones referred to in this thesis.

Einverständniserklärung

Ich bin damit einverstanden, dass meine (bestandene) Arbeit in beiden Versionen in die Bibliothek der Informatik aufgenommen und damit veröffentlicht wird.

Declaration of Consent

I agree to make both versions of my thesis (with a passing grade) accessible to the public by having them added to the library of the Computer Science Department.

Datum/Date:

Unterschrift/Signature:

SAARLAND UNIVERSITY
Software Engineering Chair
Department of Computer Science

Abstract

Network Virtualization in a Shared TCP Fast-Path by Mehrshad Lotfi Foroushani

TCP stack suffer from poor cache allocation, costly context switches, and resource sharing across different cores, making it inefficient for modern data center networking. While different techniques have been proposed to address TCP poor efficiency issues, they suffer from fundamental implications regarding maintainability, and applicability to virtualized environment, making operators and cloud tenants hesitant to deploy them.

VirtTAS is an approach to address this problem by decoupling the network stack from VMs and placing it on the hypervisor to enable streamlined TCP packet processing. By utilizing a shared TCP fast-path, tenants benefit from the multiplexing advantages, resulting in improved efficiency of CPU cycles, and decreased maintenance cost for tenants.

This thesis presents an enrichment to Virt-TAS with virtualization features by integrating Open vSwitch (OVS) to provide efficient packet processing while offering network virtualization features on top of decoupling the network stack.

We present three main contributions, including incorporating the concept of groups to improve the flexibility and scalability of Virt-TAS, leveraging OVS to take advantage of its network programmability capabilities, and optimizing the solution by using fast-path processing.

We showed that VirtTAS provides 100% increased throughput comparing to conventional Linux kernel packet processing for short messages, while consuming 40% less CPU cores, preserving precious cores for the main goal, the application logic.

CONTENT

Declaration of Authorship	iii
Abstract	v
1 Introduction	1
2 Background	5
2.1 Network Virtualization	5
2.2 Related works	6
2.2.1 OpenFlow	7
2.2.2 Open Virtual Switch (OVS)	8
2.2.3 Network Virtualization Platform (NVP)	10
2.2.4 Google Snap Platform	10
2.2.5 Virtual Filtering Platform (VFP)	11
2.2.6 TCP Acceleration as Service (TAS)	11
2.2.7 NetKernel	12
3 Design	13
3.1 Design goals	13
3.2 Design principles	14
3.3 Virt-TAS Architecture Overview	15
3.4 Slow Path	16
3.4.1 Virtual switching	16
3.4.2 Proxy	17
4 Evaluation	19
4.1 Setup	19
4.2 Single flow overhead	19
4.3 Multiplexing	20
5 Conclusion	23

List of Figures	25
List of Tables	26
Bibliography	27

CHAPTER 1

INTRODUCTION

Virtualization has become increasingly popular over the last two decades. Enabling fast provisioning, increasing the availability time, and reducing the maintenance costs for users are a few reasons of the gained popularity. Moreover, by multiplexing multiple virtual machines (VMs) on a same physical server less energy gets consumed, which makes virtualization environmental friendly and even more popular these days.

Public cloud providers such as Amazon Web Services (AWS), Microsoft Azure, and Google Cloud Platform (GCP) are expanding the range of workloads which can be spilled over virtualized environments. These vendors, in addition to providing high number of VMs with increased performance, have made it possible for public tenants to migrate their home workloads without changing the on-premises networking configurations. Furthermore, the users of these environments benefit from supplied address spaces, security groups and ACLs, scalable load balancers, bandwidth metering, QoS, and more [1].

The workloads applied by cloud tenants are bursty, which require high throughput and low latency network access. It was this requirement, which motivated cloud providers to increase their networking speeds by more than 40x and more in a span of only a few years [2]. While networking speed is growing fast, the CPU improvement has become slower, and has experienced the end of Moore's law [3]. Hence efficient packet processing is becoming more and more important.

Due to poor cache allocation, costly context switches, and resources sharing across different cores, Linux TCP packet processing is considered as inefficient for modern data center networking [4, 5]. Different techniques tried to address the issue by reducing the overheads and improving the conventional TCP stack. These techniques include bypassing kernel to enable direct NIC access from user-space [6–8] using receive side scaling (RSS) to carefully steering and processing packets on multi-cores architecture [4, 9], and offloading packet processing to NICs with computation capabilities [2, 5, 10, 11].

Although these techniques have improved the performance of end-host packet processing, they suffer from two fundamental implications when it comes to providing high throughput and low latency network access for applications residing on VMs.

Implication I. Majority of these techniques are proposed only for applications running on bare metal servers, and they do not provide the rich virtualization features required in multi-tenant datacenters. Hence, cloud providers hesitate to deploy the techniques for their tenants and applications still use the conventional TCP packet processing.

Implication II. Cloud tenants are only concerned with performance of their applications, and they do not want the maintenance of the network stack to fall out on their shoulders. Deploying new packet processing stacks requires extensive testing and comes with the risk of misconfigurations. Therefore, cloud tenants hesitate to deploy the new techniques in trade with high throughput and low latency.

Goal In this thesis, we ask the following question: *How can we provide virtualization features on top of the modern stacks?* In response, we enrich Virt-TAS with virtualization features. Virt-TAS is a TCP acceleration service (TAS) for virtualized environments targeted at applications that require low latency and high throughput. It runs as a service alongside the host and provides a fast-path for common send and receive operations, through which it reduces the virtualization overheads incurred by the hypervisor and guest operating system [12].

To avoid reinventing the wheel, instead of implementing a virtual switch from scratch, we complement Virt-TAS by integrating Open vSwitch (OVS) [13]. OVS is a virtual switch that is deployed in many data center networks. It is also a component of VMware's NSX product, used by thousands of customers [14].

NetKernel, is another system that provides network stack as a service in the cloud [15, 16]. Like Virt-TAS, Netkernel decouples the network stack from the guest OS, (i) to simplify deployment and maintenance for tenants, and (ii) to increase efficiency of TCP packet processing for cloud tenants. Netkernel's goal is to show the feasibility of decoupling network stack from VMs, whereas Virt-TAS focuses on providing efficient packet processing while providing network virtualization features on top of decoupling network stack.

This thesis makes the following contributions.

- We have made an extension to the memory layout of Virt-TAS to incorporate the concept of *groups*. Inspired by Amazon EC2's target groups [17], a group in Virt-TAS is a set of virtual machines (VMs) that trust each other and can share memory. In this design, each VM in Virt-TAS is assigned to a group, and a group

can be assigned to multiple VMs. This extension improves the flexibility and scalability of the Virt-TAS system and allows for more efficient memory sharing among VMs.

- We integrate OVS to a TCP Acceleration Service (TAS) to leverage the network programmability features provided by the OpenFlow interface. In the new design, each VM is connected to OVS through a shared memory, managed by TAS. TAS, on the other hand, facilitates communication between the NIC and OVS, as well as between applications on VMs and the hypervisor, through lock-free queues on shared memories. The design enables us to take advantage of OVS's network programmability capabilities and benefit from improved network performance through streamlining TCP packet processing enabled by TAS.
- As a further optimization, we leverage fast-path processing in our solution. Virt-TAS processes the first packet of each connection in Open vSwitch (OVS) and stores the connection information for the flow in the cache of TAS for bridge-only connections. The subsequent packets utilize the fast-path processing. This optimization results in improved performance and reduced overhead in the processing of network flows.

The remainder of this thesis is structured as follows: Chapter 2 provides a comprehensive background on network virtualization, including a review of the state-of-the-art solutions. Chapter 3 delves into the design of Virt-TAS, as well as outlining the design requirements. Chapter 4 explains the implementation of Virt-TAS. In chapter 5, we evaluate Virt-TAS based on three predefined requirements. Finally, chapter 6 provides a conclusion of our work, summarizing our findings and contributions to the field of network virtualization.

CHAPTER 2

BACKGROUND

2.1 Network Virtualization

Network virtualization has gained an extreme importance in data center networking as it enables cloud providers to enhance the flexibility, scalability and resource utilization of the network infrastructure. With network virtualization multiple virtual networks can be deployed on a single physical infrastructure. In this section, we explain what virtualized network environments are and how they are different from the conventional non-virtualized environments.

In a non-virtualized network environment, networking components are directly tied to the underlying physical hardware. Each networking component, such as a switch or

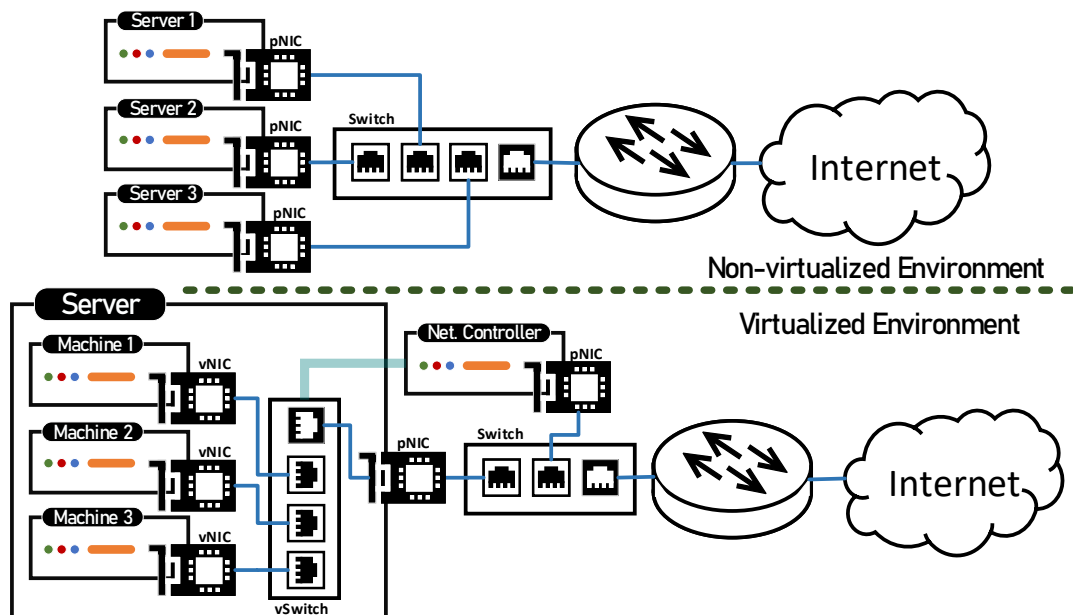


FIGURE 2.1: Non-virtualized vs. virtualized network environments

router, corresponds to a specific physical device. Therefore, the network as a whole is limited by the available hardware resources. Furthermore, configuration in these environments is a hard and time-consuming task. Each networking vendor has its own interface, and configuring each device requires the corresponding expertise. In addition to that, adding or removing a machine requires multiple configurations to be set up in a box-by-box manner, which introduces an excessive operation overhead and increases the risk of misconfiguration [9, 18].

Network virtualization aims to solve these issues by separating and abstracting network infrastructure from the underlying hardware. Through a virtualization layer, users can create virtual networks, with arbitrary service models, topologies, and addressing scheme, on top of the same physical network devices. Furthermore, a global abstraction enables management and configuration of these virtual networks [19].

Figure 2.1 illustrates differences between a non-virtualized and a virtualized networking environment. In a non-virtualized setup, the kernel directly routes packets from applications to a physical Network Interface Card (pNIC), binding the applications to the physical networking components. Whereas, in a virtualized setup, software-based virtual NICs (vNICs) and virtual switches (vSwitches) are used to abstract and separate packets from the underlying hardware through an encapsulation layer. Moreover, a network controller in collaboration with virtual switches configures policies for encapsulation and forwarding decisions in virtualized environments.

2.2 Related works

To provide network virtualization for cloud tenants, public cloud providers, such as Amazon Web Services (AWS), Microsoft Azure, and Google Cloud Platform (GCP), have taken different approaches. While Microsoft has been following a hardware-assisted procedure to provision networking demands [1, 2], Google has been utilizing a software-only approach [9]. In this section, we first introduce two seminal open-source projects in network virtualization, namely OpenFlow (Section 2.2.1), and Open Virtual Switch (Section 2.2.2). We then elaborate on the state-of-the art network virtualization used by VMWare (Section 2.2.3), Google (Section 2.2.4), and Microsoft (Section 2.2.5). Finally, we explore two recent academic project: TCP Acceleration as an OS Service (TAS) (Section 2.2.6), and NetKernel (Section 2.2.7).

Table 2.1 presents a comparison of the different approaches, highlighting their level of support and benefits across the following four categories: (i) Native virtualization support, (ii) Multiplexing of workloads from different VMs, (iii) capability to support

	Virtualization support	Multiplexing benefits	Hardware offload	TCP acceleration on hypervisor
NVP	✓	✗	✗	✗
Snap	✓	✗	✗	✗
VFP	✓	✗	✓	✗
TAS	✗	✗	FlexTOE[5]	✓
NetKernel	✓	✓	✗	✗
VirtTAS	✓	✓	FlexTOE[5]	✓

TABLE 2.1: Comparison of different network stack architectures

hardware offloads, and (iv) enhanced efficiency of packet processing by leveraging acceleration service.

As shown in the table, VirtTAS offers native virtualization support, as well as the ability to multiplex workloads among different virtual machines and process packets efficiently by utilizing acceleration services. Moreover, based on a recent research on offloading TCP acceleration service to hardware programmable NICs[5], it is feasible to further enhance VirtTAS and take advantages of hardware support.

2.2.1 OpenFlow

OpenFlow was first proposed to enhance the implementation of new protocols on top of the production networks. It consists of a standardized interface to add or remove entries from flow tables residing on the Ethernet switches [19].

Flow tables are the fundamental data structure of an OpenFlow switch. They consist of flow entries and proper actions to the packets of each flow. A flow is a set of packets that have similar characteristics. It could be a TCP connection, all packets from a particular port number, or packets with the same VLAN tag. An OpenFlow switch matches the packets with the specified flows and applies the action of that flow. For example, all packets from a particular sender could be dropped by the OpenFlow switch.

The three basic actions that are supported by an OpenFlow switch consists of (i) forwarding packets to a given port (or ports), (ii) encapsulating and forwarding packets to a controller, and (iii) dropping packets.

In summary, an OpenFlow switch has the following three main components:

1. *A Flow-table residing on the Ethernet switch:* This flow table maintains information about existing flows and actions required to be applied to the corresponding packets.
2. *Secure interface between a controller and the switch:* The interface is used when a packet has no corresponding entry in the flow table. In such scenarios, the packets will be forwarded through this interface to a controller. The controller can then decide for the corresponding action.
3. *Standardized OpenFlow interface to program flow tables:* With this standardized interface, a controller can update the flow-tables on Ethernet switches. Therefore, not all the packets are required to be transferred to the controller, hence higher throughput can be achieved.

Further detailed specifications of an OpenFlow switch are defined by *Open Networking Foundation* [20].

VirtTAS supports OpenFlow interface to enable network programmability at end-host.

2.2.2 Open Virtual Switch (OVS)

Conventional virtual switches are concerned only with providing a basic network connection for VMs through L2 networks. This was changed with the advent of network virtualization. Today's virtual switches provide most network services for VMs, and the physical networks only transmit IP-tunneled packets. Therefore, VMs are no longer tightly connected to physical datacenter networks and the workloads benefit from higher scalability and mobility.

As the complexity of network virtualization increased, the need for a production-level multi-layer virtual switch became a necessity. Open vSwitch (OVS) was designed and implemented to address this need. It enables network virtualization through OpenFlow and supports standard interface management protocols such as NetFlow[21], sFlow[22].

Open vSwitch has a long history of development and in order to keep up with the production workloads and increase the development velocity, the following key design choices were made by OVS community:

1. **Use tuple search classification..** Tuple search classification supports constant time updates. Constant time updates are essential in virtualized environments, where a controller sends flow updates multiple times a second. Furthermore, it

consumes memory linear in the number of flows, and it can be generalized to all combinations of packet header fields.

2. **Split implementation between Kernel and userspace.** Kernel-based networking makes the development and release cycles lengthy. Moreover, OVS would not be accepted as a contribution to upstream Linux, if it was implemented solely in Kernel. Thereby, OVS uses kernel as a match and action cache and the main classification labor is done in userspace by a virtual switch daemon.
3. **Benefit from caching.** OVS uses a two-layers caching mechanism to match packets in the Kernel. A Micro-flow cache layer matches on specific flows with exact headers fields, and a Mega-flow cache layer matches on a wildcard of flows. These wildcards are generated by the cross-product of OpenFlow match and action tables.

Although OVS has a positive impact on making virtual switches available for the networking community, there are a few shortcomings which makes it not suitable for large scale deployments.

1. **Explicit tunnel interface.** OVS uses an extra interface to setup tunneling actions, and it hardcodes models of tunneling to dataplane, rather than enabling a controller to specify details of tunneling. Therefore, adding complicated tunneling logics, such as ECMP routing, is difficult with OVS.
2. **No support for multi-controller model.** In a multi-controller software-defined network, a virtual switch transmits in-bound packets in the reverse direction of out-bound packets to make the states consistent, whereas OVS only supports forward transmission.
3. **No support for stateful actions** Actions are not stateful, there are trackers supported by OVS but it does not support stateful actions.
4. **Performance** Packets that are not matched in the fast-path cache are handled by the slowpath and in specific cases, these packets are forwarded to a controller. This makes OVS vulnerable to DDoS attacks and cache pollution.

Despite the limitations of OVS, we decided to integrate OVS into our shared TCP fast-path, and take the benefits of its long history of development and deployment. Addressing the limitations of OVS is beyond the scope of this thesis and we believe the OVS community would address these issues in the future works.

2.2.3 Network Virtualization Platform (NVP)

Network Virtualization Platform (NVP) is one of the first commercial virtualization platform based on SDN that appeared with the market demand for network virtualization and research on Software Defined Networking (SDN) [23]. NVP was proposed to address virtualization primitives such as VLAN (Virtualized L2 domain), VRFs (Virtualized L3 FIB), NAT (Virtualized IP address space), and MPLS (Virtualized Path) needed by VMware's users.

NVP allows users to setup, and configure different virtual networks, each with different addressing space, independent services model, and topologies on top of the same physical network. It hides the topology of the underlying physical network from the users, and specific aspects of the forwarding devices becomes unknown for the users. In NVP, a hypervisor layer translates each tenant's network configurations to low-level instruction set, which can be installed by virtual switches. NVP deploys OVS as its virtual switch. The packets are tunneled by virtual switches over the physical network, and the physical switch only transmit tunneled packets between hypervisors and gateways.

Virt-TAS like NVP, enables virtualization primitivs through use of OVS at end-host. NVP, focuses on providing network connectivity to VMs and network programmability for operators, whereas Virt-TAS aims to supply VMs with high throughput and low latency by streamlining packet processing as a service on hypervisor in addition to providing the same features.

2.2.4 Google Snap Platform

As one of the primary cloud providers with internet-based products that are used by a wide range of users, Google has prioritized the following three main requirements for host networking in their data centers.

1. **Fast development and release cycles.** The network bandwidth is continuously increasing and delivering novel approaches for edge switching and bandwidth management is inevitable. Thus, a system that enables fast delivery of innovative paradigms is required.
2. **Virtualization features.** Rich virtualization features must be provided on top of host networking to enable cloud computing.
3. **Low latency and high throughput.** Distributed data-intensive applications require low latency and high throughput for their communication. As a consequence, the host networking system should be optimized with regard to these applications.

To reach the requirements of data-center networking and to mitigate the drawbacks of using the kernel in the development of end-host network stack, Google has proposed Snap [9]. Snap is a userspace networking system, which was initially inspired by microkernel architecture. By residing in userspace, the development of new features has become faster and upgrades to the networking stack can be applied transparently. Moreover, in comparison to the Linux Kernel networking stack, it provides higher throughput.

2.2.5 Virtual Filtering Platform (VFP)

Microsoft Azure has deployed Virtual Filtering Platform (VFP) as its programmable virtual switch on more than 1M hosts to provide virtualization features [1]. VFP extends Microsoft Hyper-V's switch [24] with the switch logic, and it implements match and actions tables with multiple layers to support multi-controller programmability. VFP takes advantage of a new programming model, which unlike OpenFlow consists of complex actions.

The key advantage of VFP is its performance which is enabled by using a just in time compiler, metadata parsing and touching packet once all decisions are made.

Unfortunately the implementation details of VFP is not publicly available and we cannot evaluate performance benefits of extending VFP to streamline packet processing on endhost. Nevertheless, we believe the main idea is still applicable to VFP. We believe by providing packet processing as a service at hypervisor, the tenants do not need to provision higher than their demands.

2.2.6 TCP Acceleration as Service (TAS)

TCP acceleration techniques are used in data center networking to increase TCP's performance on high-bandwidth, low-latency networks. Some of these techniques are as follows:

- **TCP/IP Offload Engine (TOE).** TOEs are hardware-based solutions that offloads TCP/IP processing from the host CPU to a dedicated hardware on network interface card (NIC), therefore lowering CPU usage and accelerating network performance [5, 25–27].
- **Kernel bypass.** Kernel bypass is another technique used to directly access the NIC from the application in userspace. This technique enables low latency communication between the application and the NIC, as it eliminates the need for the data to pass through the kernel's network stack [4, 9, 28].

- **Congestion Control.** Different congestion control are used to manage network traffic and to prevent congestion and data loss [29, 30].

TAS is an acceleration service, which splits TCP packet processing into two components: a light weight fast-path optimized for common case client-server RPCs, and a heavier stack slow-path which drives congestion control, connection setup.

The fast-path benefits from different streamlining opportunity. The user can scale the dedicated cores for efficient and fast processing of packets.

2.2.7 NetKernel

Netkernel is another system that offers network stack as a service in virtualized environments. Similar to VirtTAS, Netkernel decouples the network stack from the VMs. Thus, (i) tenants benefit from fast deployment and maintenance of the network stack, and (ii) cloud providers are able improve the efficiency of TCP stack processing stacks transparently.

In the design of Netkernel, the BSD socket APIs are redirected to a full NetKernel socket implementation through a library called *GuestLib*. The library is deployed as a kernel patch, and it is the only modification made to the users' virtual machines. Network stacks are set up by the operator on different virtual machines inside the same server as *Network Stack Modules (NSMs)*. The NSMs are connected from one side to users' virtual machines through inter-VM shared-memory and from the other side to a virtual switch through conventional tap devices.

The design of NetKernel is limited by processing network stack inside virtual machines. Furthermore, the full performance benefits claimed by NetKernel is only achieved when cloud providers deploy Single Root I/O virtualization (SR-IOV) in their datacenter. SR-IOV ties applications to the underlying hardware, and is not favorable to all cloud providers. Whereas, VirtTAS focuses on improving the performance of packet processing for application residing on user virtual machines through streamlining the packet processing on the hypervisor. VirtTAS benefits from the multi-core achitecture and uses kernel-bypass to achieve line-rate performance.

CHAPTER 3

DESIGN

3.1 Design goals

Virt-TAS is designed to fulfill the following goals:

1. **Low latency and high throughput** Modern data-centers host distributed data-intensive applications, that require low latency and high throughput. Workloads of these applications consist of short-lived flows with stringent low latency requirements, as well as long-lived flows with high throughput needs. VirtTAS should provide the latency and throughput needs of these applications, and it should enable cloud providers to offer and maintain Service Level Agreements (SLAs).
2. **Compatibility** The networking interfaces should be unmodified so that applications residing on the VMs could benefit from VirtTAS without any modification. Moreover, to make VirtTAS programmable through existing state-of-the-art network controllers, it should support OpenFlow API.
3. **Coexistence with traditional Kernel networking** The new networking framework should be capable of working in conjunction with the traditional virtual machine networking system. This allows tenants to execute certain workloads using the available kernel networking, while utilizing VirtTAS for performance-critical workloads.

In order to ensure seamless coexistence, it is necessary for VirtTAS to operate without adversely impacting the performance of other applications that are running on the same machine using kernel networking.

4. **Resource efficiency** As VirtTAS shares fate with the hypervisor, resource conversation is highly critical. VirtTAS should maximize resource efficiency to maintain precious resources for the main goal of the hypervisor, running user workloads.

5. **Scalability** Scalability becomes a challenging problem when different tenants with different network topologies and networking demands reside on the same hypervisor. VirtTAS should keep up with the increasing number of flows. It should also support an increasing number of flow updates from a controller.
6. **Packets confidentiality** Applications residing on one VM must be prevented to access packets from other VMs. VirtTAS should provide this protection of packets which may contain sensitive information from unauthorized access, disclosure, or interception.

In order to advance the VirtTAS framework, several other goals have been identified for future work. The pursuit of the following goals is expected to enhance the functionality and effectiveness of our framework.

- **Mobility** VMs could be migrated to different locations, thus VirtTAS should allow VMs to continue communicating over the network, despite their migration. VirtTAS should completely decouple applications from physical networking infrastructures at their location.
- **Isolation** Applications of one tenant must be prevented to violate other tenants' SLAs. Thus, VirtTAS must provide an isolation mechanism that ensures this isolation exist.

3.2 Design principles

To achieve the aforementioned goals, we benefit from a set of design principles to provide network virtualization features in a shared TCP fastpath.

- **Fast-path/Slow-path separation** We integrated TCP Acceleration Service into our design which is an optimized and streamlined processing of TCP network packets that follow a commonly used code path. As seen in the related work, fast path processing is a common design approach for network applications to quickly process packets that do not require complex processing or special handling. By separating these packets from the ones that require special treatment, fast path processing enables the network to handle a large number of packets efficiently and with minimal overhead.
- **Proxy design principle** A pair of proxies runs on the host and guest VM and acts as an intermediary between the VMs and the TCP acceleration service. When

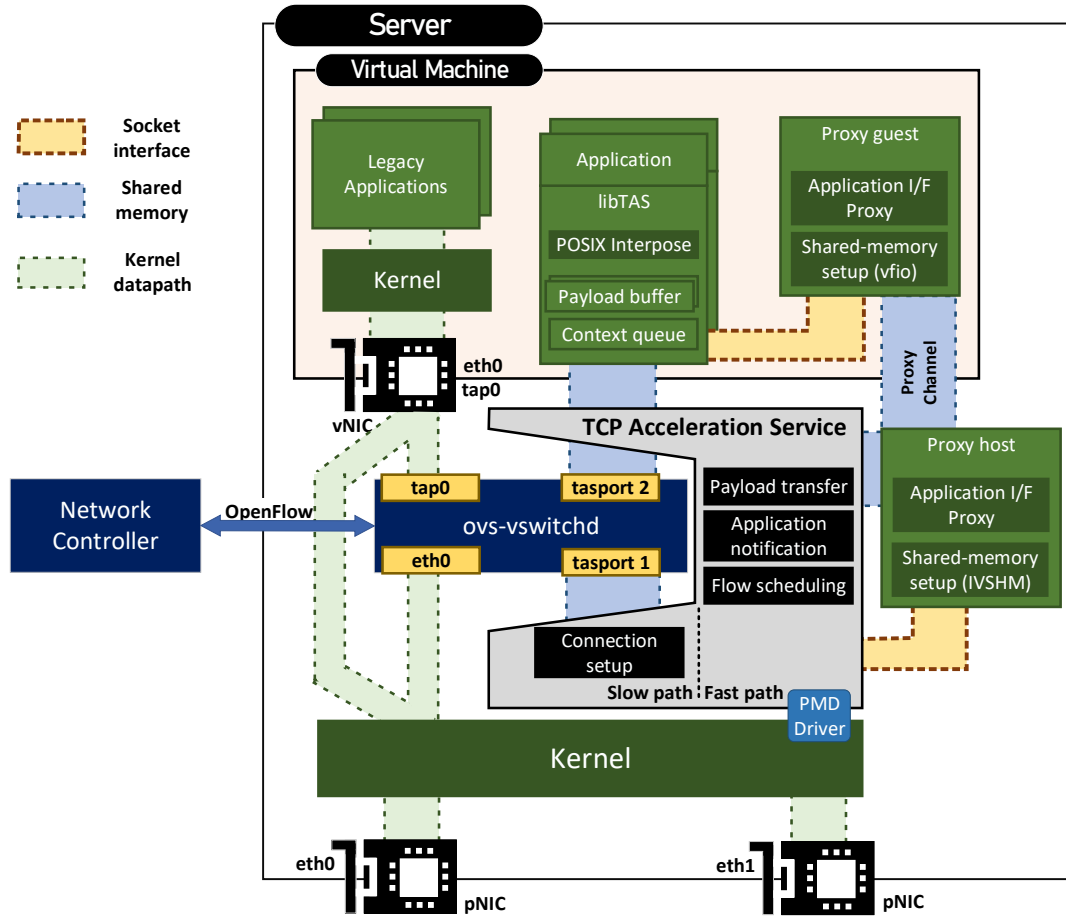


FIGURE 3.1: Virt-TAS design overview

a new application wants to connect to TAS, it sends a context request, which is intercepted by the proxy. The proxy forwards the request to TAS and returns the response to the VM.

Using proxy enables VirtTAS to be integratable to virtual manager programs such as virt manager [31], or Virtual Box [32] as it acts as a gatekeeper, controlling VM access to TAS and providing an additional layer of security.

3.3 Virt-TAS Architecture Overview

Similar to TAS, Virt-TAS devides the network stack into three different components: (i) A fast-path running on hypervisor, (ii) a slow-path that has components both on the hypervisor and the VMs, and (iii) an application library, which runs entirely on the VM side.

Figure 3.1 depicts the overview architecture of Virt-TAS. The packets are first received by the fast-path component from a physical NIC or through shared memory from applications running on VMs. The fast-path either has received instruction from OVS on how to handle the received packet based on its flow or it has not. In the first scenario, the fast-path applies the action provided by OVS. The actions ranges from dropping the packet to forwarding it to specific application connected to Virt-TAS or to the physical NIC. In the second scenario, where fast-path has no information about the flow of the received packet, the packet is delivered to OVS through lock-free shared memory. OVS can then decide how to treat packets. On the other side, the application interface has the responsibility to initialize the connection to Virt-TAS and implement POSIX network socket API for applications.

By integrating OVS, we can control and program the forwarding dataplane through OpenFlow protocol. OpenFlow enables network operators to add, remove, update entries, and to monitor statistics on flow tables. OVS takes OpenFlow tables from an SDN controller, matches the received packets to these flow tables, and applies all actions needed to be taken. The outcome is then cached by the fast-path component of Virt-TAS. This significantly simplifies the fast-path, as it allows the fast-path to be agnostic to OpenFlow specifics. On the other side, from the perspective of an OpenFlow controller, processing packets in Virt-TAS is an invisible implementation detail. In its perspective, all packets are matched against multiple flow tables, and the corresponding actions are applied among entries that satisfy all conditions and have the highest priority.

A userspace library redirects the POSIX socket API transparently, so that applications remain unmodified. Through this library the sockets calls are transmitted to the hypervisor instead of being transfered to the VMs network stack. Furthermore, By rewriting applications to use low level API of libTAS higher performance can be achieved.

3.4 Slow Path

The slow-path in VirtTAS operates on a dedicated thread. It is responsible for handling exceptional cases related to policy decisions and virtualization management, as well as congestion control, connection control, and connection setup and teardown.

3.4.1 Virtual switching

VirtTAS uses OVS as its virtual switch to direct traffic to the appropriate VM. The first packet of a flow is not recognised by the fast-path, therefore the responsibility of the

packet is handed over the slow-path and OVS by the fast-path. This forwarding is done via lock-free shared memories. OVS then updates the corresponding flow table with the flows information and communicates the decisions to VirtTAS. Subsequent packets for the same flow are not rerouted to the slow-path and the decisions are available in the cache of fast-path.

3.4.2 Proxy

The Proxy Host is a service that operates on the hypervisor as an intermediary between virtual machines and the TAS. It uses the libtas library to connect to the running TAS instance on the hypervisor. The Proxy Host determines the number of supported groups from the TAS instance and sets up a UNIX socket for each group. The path of the created socket for the desired group is then used to launch virtual machines with fast-path support.

```
1 qemu-system-x86_64 \  
2 ...  
3 -chardev socket,path="PROXY_SOCKET_PATH",id="tas" \  
4 -device ivshmem-doorbell,vectors=1,chardev="tas" \  
5 ...
```

LISTING 3.1: QEMU system x86-64 command with ivshmem-doorbell and socket interface for connection to proxy host

runs a request to a server, it is sent to the proxy server first. The proxy server then makes a request to the server on behalf of the client, and returns the response back to the client.

CHAPTER 4

EVALUATION

4.1 Setup

Our evaluation cluster comprises two Intel(R) Xeon(R) Gold systems, each with 48 cores and 48 hyperthreads running at 2.40GHz. The systems are equipped with 260 GB RAM, 72 MB aggregate cache, and an Intel E810 100Gb Ethernet adapter. One of the systems serves as the server, while the other is used for clients. All the machines run Debian GNU/Linux 11 (bullseye) with kernel version 5.15, and we use Open vSwitch version 3.1.0. We build TAS using DPDK 21.11.

4.2 Single flow overhead

Throughput In this experiment, we aimed to determine the overhead of VirtTAS compared to running an application in a bare-metal approach without virtualization. To achieve this, we conducted a Micro-RPC benchmark using three different approaches: 1. baremetal TAS, 2. VirtTAS without using groups, and 3. VirtTAS with 16 groups.

Figure 4.1 shows the throughput of the client each with 8 threads. As illustrated in the figure, the throughput for Group VirtTAS is less than 10% for packet sizes below 256 bytes. This indicates that VirtTAS incurs reasonable overhead when compared to running the application in a bare-metal environment. In other words, the results also suggest use of VirtTAS has comparable performance to non-virtualized environment in scenarios with smaller packet sizes.

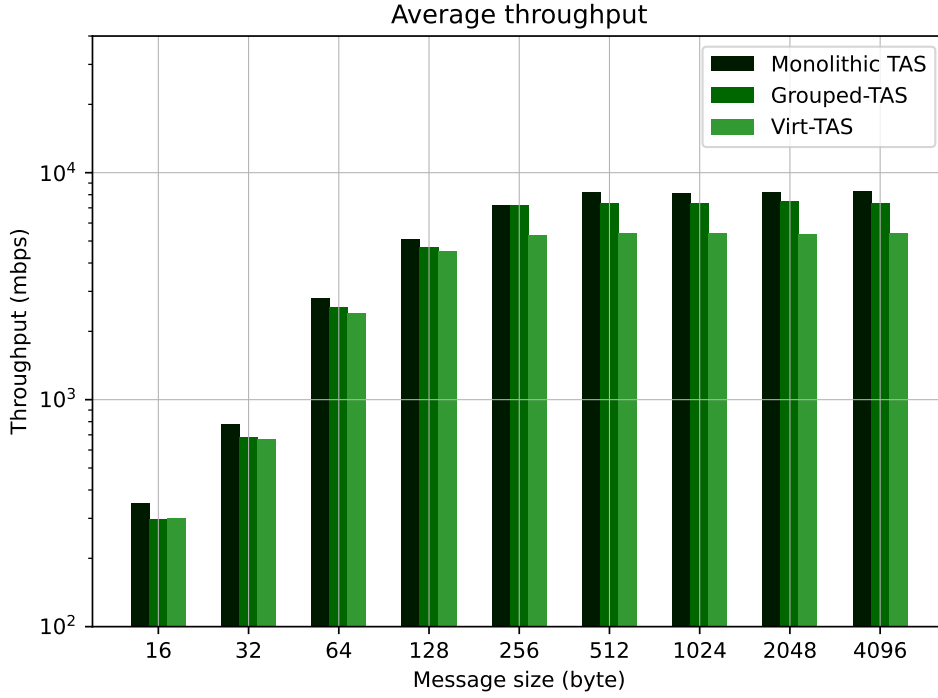


FIGURE 4.1: Single context overhead.

4.3 Multiplexing

The goal of multiplexing experiment is to validate the enhancement of VirtTAS in terms of higher throughput and resource efficiency. Specifically, the experiment investigates the ability of VirtTAS to benefit from multiplexing of packet processing across multiple VMs through a shared TCP fast-path.

To this end, we conduct three different scenarios, as illustrated in Figure 4.2. In each of the scenarios, we configure four VMs across two servers. The servers in this experiment are interconnected directly. Subsequently, we execute Micro-RPC benchmarks between four pairs of VMs, with message sizes from 16 to 4096.

In the first scenario (*kernel*), we follow the conventional setup where applications running on the VMs utilize the conventional kernel networking. In this scenario, the hypervisor establishes connectivity to VMs through TAP devices, while a virtual switch applies virtualization policies to the network. For the setup, we allocate three cores to each VM. The server and client applications are using three threads to leverage all available cores. Additionally, the virtual switch runs on a separate core. Overall, the scenario employs a total of 13 cores.

In the second scenario (*Independent-TAS*), we execute the TCP acceleration service on each VM, utilizing two cores for the service and one core for the application itself. The goal is to optimize TCP packet processing within the VMs, to increase throughput and reduce latency for the applications. Similar to the first scenario, the virtual switch applies virtualization policies to packets on the hypervisor, and the experiment consumes 13 cores.

In the final scenario (*Virt-TAS*), we deployed VirtTAS on the hypervisor, dedicating four cores exclusively to its operations. In this case, each VM is allocated a single core for application processing. Notably, this scenario requires only 8 cores, which is nearly 40% less cores comparing to previous scenarios.

Figure 4.3 illustrates the average throughput of Micro-RPC benchmark in three aforementioned different scenarios. The analysis indicates that, for message sizes smaller than 256, VirtTAS outperforms the alternative methods. Specifically, for message of size 256, VirtTAS yields a 25% higher average throughput than the kernel case, and for message size 64, VirtTAS achieves $2\times$ the throughput of the kernel case.

Based on the results presented in Figure 4.3, it is evident that in scenarios where the cloud provider does not offer the TCP Acceleration Service (TAS), they can still take advantage of fast-path TCP by deploying TAS on their VMs for workloads with message sizes of less than 32 bytes.

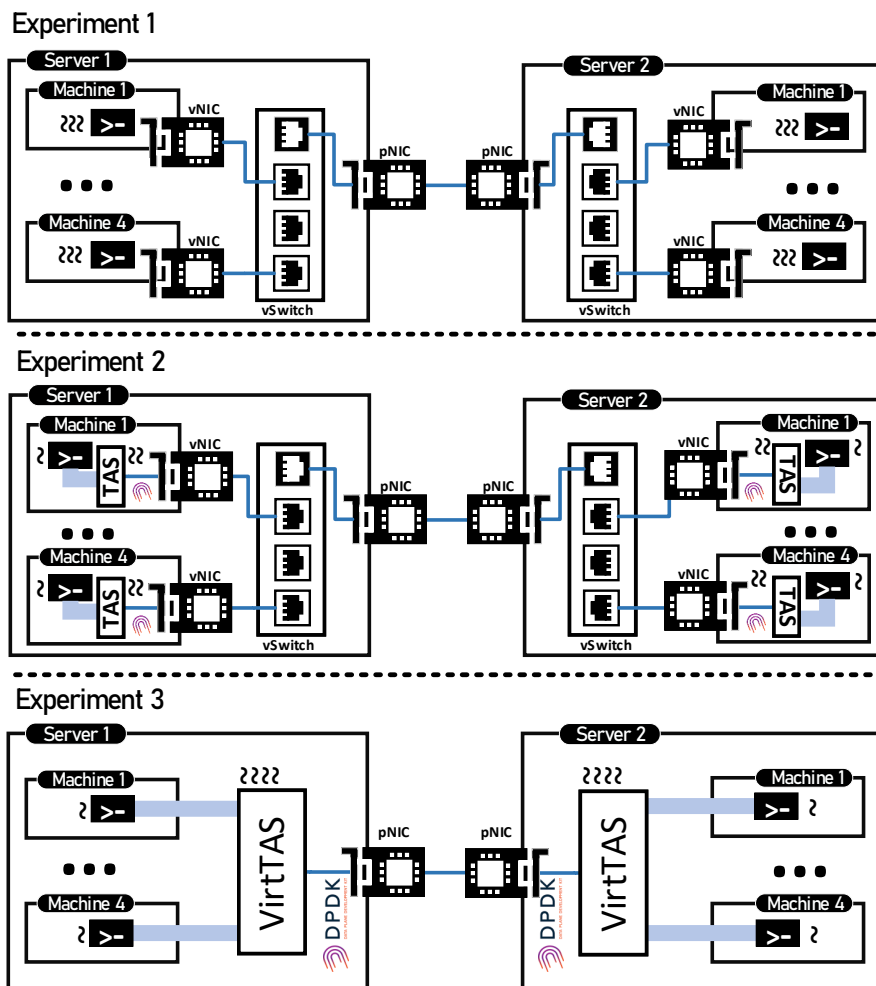


FIGURE 4.2: Multiplexing experiment setup in three different scenarios (i) kernel, (ii) Independent-TAS, and (iii) Virt-TAS.

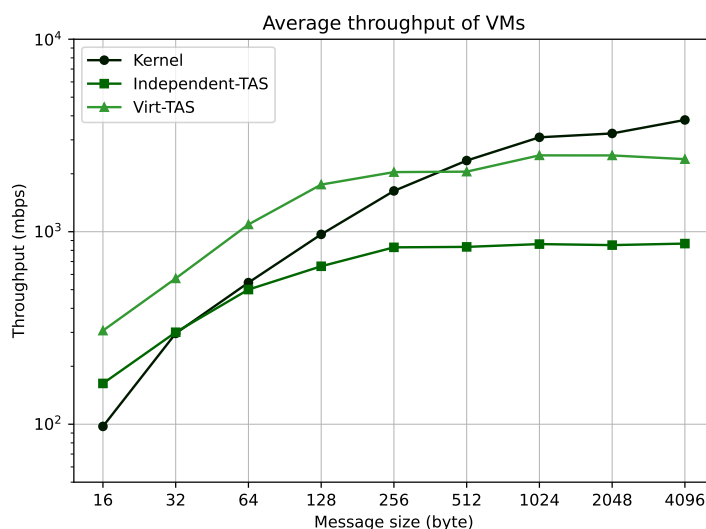


FIGURE 4.3: Average throughput of Micro-RPC benchmark in three different scenarios (i) kernel, (ii) Independent-TAS, and (iii) Virt-TAS

CHAPTER 5

CONCLUSION

In conclusion, VirtTAS provides a solution to the problem of inefficient TCP packet processing by decoupling the network stack from virtual machines and placing it on the hypervisor. By utilizing a shared TCP fast-path, VirtTAS offers multiplexing advantages, improved efficiency of CPU cycles, and decreased maintenance costs for tenants.

We extended VirtTAS by integrating Open vSwitch (OVS) to provide efficient packet processing and network virtualization features. We improved the flexibility and scalability of TAS by introducing groups. Furthermore, we leveraged OVS to improve network programmability capabilities, and we optimized the solution by using fast-path processing.

The results show that VirtTAS provides a 100% increase in throughput compared to conventional Linux kernel packet processing for short messages, while consuming 40% less CPU cores. This ensures that valuable CPU resources can be allocated to the main application logic. Overall, VirtTAS with OVS integration enables efficient and flexible TCP packet processing in virtualized environments.

LIST OF FIGURES

Figure 2.1	Non-virtualized vs. virtualized network environments	5
Figure 3.1	Virt-TAS design overview	15
Figure 4.1	Single context overhead.	20
Figure 4.2	Multiplexing experiment setup in three different scenarios (<i>i</i>) kernel, (<i>ii</i>) Independent-TAS, and (<i>iii</i>) Virt-TAS.	22
Figure 4.3	Average throughput of Micro-RPC benchmark in three different scenarios (<i>i</i>) kernel, (<i>ii</i>) Independent-TAS, and (<i>iii</i>) Virt-TAS	22

LIST OF TABLES

Table 2.1	Comparison of different network stack architectures	7
-----------	---	---

BIBLIOGRAPHY

- [1] Daniel Firestone. {VFP}: A virtual switch platform for host {SDN} in the public cloud. In *14th USENIX Symposium on Networked Systems Design and Implementation (NSDI 17)*, pages 315–328, 2017.
- [2] Daniel Firestone, Andrew Putnam, Sambhrama Mundkur, Derek Chiou, Alireza Dabagh, Mike Andrewartha, Hari Angepat, Vivek Bhanu, Adrian Caulfield, Eric Chung, et al. Azure accelerated networking:{SmartNICs} in the public cloud. In *15th USENIX Symposium on Networked Systems Design and Implementation (NSDI 18)*, pages 51–66, 2018.
- [3] Hadi Esmaeilzadeh, Emily Blem, Renee St. Amant, Karthikeyan Sankaralingam, and Doug Burger. Dark silicon and the end of multicore scaling. In *Proceedings of the 38th annual international symposium on Computer architecture*, pages 365–376, 2011.
- [4] Antoine Kaufmann, Tim Stamler, Simon Peter, Naveen Kr Sharma, Arvind Krishnamurthy, and Thomas Anderson. Tas: Tcp acceleration as an os service. In *Proceedings of the Fourteenth EuroSys Conference 2019*, pages 1–16, 2019.
- [5] Rajath Shashidhara, Tim Stamler, Antoine Kaufmann, and Simon Peter. {FlexTOE}: Flexible {TCP} offload with {Fine-Grained} parallelism. In *19th USENIX Symposium on Networked Systems Design and Implementation (NSDI 22)*, pages 87–102, 2022.
- [6] Adam Belay, George Prekas, Ana Klimovic, Samuel Grossman, Christos Kozyrakis, and Edouard Bugnion. {IX}: a protected dataplane operating system for high throughput and low latency. In *11th USENIX Symposium on Operating Systems Design and Implementation (OSDI 14)*, pages 49–65, 2014.
- [7] EunYoung Jeong, Shinae Wood, Muhammad Jamshed, Haewon Jeong, Sunghwan Ihm, Dongsu Han, and KyoungSoo Park. {mTCP}: a highly scalable user-level {TCP} stack for multicore systems. In *11th USENIX Symposium on Networked Systems Design and Implementation (NSDI 14)*, pages 489–502, 2014.
- [8] George Prekas, Marios Kogias, and Edouard Bugnion. Zygus: Achieving low tail latency for microsecond-scale networked tasks. In *Proceedings of the 26th Symposium on Operating Systems Principles*, pages 325–341, 2017.

- [9] Michael Marty, Marc de Kruijf, Jacob Adriaens, Christopher Alfeld, Sean Bauer, Carlo Contavalli, Michael Dalton, Nandita Dukkupati, William C Evans, Steve Gribble, et al. Snap: A microkernel approach to host networking. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles*, pages 399–413, 2019.
- [10] Mina Tahmasbi Arashloo, Alexey Lavrov, Manya Ghobadi, Jennifer Rexford, David Walker, and David Wentzlaff. Enabling programmable transport protocols in {High-Speed}{NICs}. In *17th USENIX Symposium on Networked Systems Design and Implementation (NSDI 20)*, pages 93–109, 2020.
- [11] Jiaxin Lin, Kiran Patel, Brent E Stephens, Anirudh Sivaraman, and Aditya Akella. {PANIC}: A {High-Performance} programmable {NIC} for multi-tenant networks. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*, pages 243–259, 2020.
- [12] Florian Bauckholt. virt-tas: Virtualization capabilities for the tcp acceleration service, 2021.
- [13] Ben Pfaff, Justin Pettit, Teemu Koponen, Ethan Jackson, Andy Zhou, Jarno Rajahalme, Jesse Gross, Alex Wang, Joe Stringer, Pravin Shelar, et al. The design and implementation of open {vSwitch}. In *12th USENIX symposium on networked systems design and implementation (NSDI 15)*, pages 117–130, 2015.
- [14] William Tu, Yi-Hung Wei, Gianni Antichi, and Ben Pfaff. Revisiting the open vswitch dataplane ten years later. In *Proceedings of the 2021 ACM SIGCOMM 2021 Conference*, pages 245–257, 2021.
- [15] Zhixiong Niu, Qiang Su, Peng Cheng, Yongqiang Xiong, Dongsu Han, Keith Winstein, Chun Jason Xue, and Hong Xu. Netkernel: Making network stack part of the virtualized infrastructure. *IEEE/ACM Transactions on Networking*, 2021.
- [16] Zhixiong Niu, Hong Xu, Dongsu Han, Peng Cheng, Yongqiang Xiong, Guo Chen, and Keith Winstein. Network stack as a service in the cloud. In *Proceedings of the 16th ACM Workshop on Hot Topics in Networks*, pages 65–71, 2017.
- [17] Target groups for your application load balancers. <https://docs.aws.amazon.com/elasticloadbalancing/latest/application/load-balancer-target-groups.html>. Accessed: 2023-01.
- [18] DW Cearley, D Scott, J Skorupa, and TJ Bittman. Top 10 technology trends, 2013: cloud computing and hybrid it drive future it models. *Gartner, February*, 2013.
- [19] 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.
- [20] OpenFlow Switch Specification. Version 1.0. 0 (wire protocol 0x01). *Open Networking Foundation*, 2009.
- [21] Benoit Claise. Cisco systems netflow services export version 9. Technical report, 2004.
- [22] Mea Wang, Baochun Li, and Zongpeng Li. sflow: Towards resource-efficient and agile service federation in service overlay networks. In *24th International Conference on Distributed Computing Systems, 2004. Proceedings.*, pages 628–635. IEEE, 2004.
- [23] Teemu Koponen, Keith Amidon, Peter Baland, Martín Casado, Anupam Chanda, Bryan Fulton, Igor Ganichev, Jesse Gross, Paul Ingram, Ethan Jackson, et al. Network virtualization in multi-tenant datacenters. In *11th USENIX Symposium on Networked Systems Design and Implementation (NSDI 14)*, pages 203–216, 2014.
- [24] Hyper-v virtual switch. [https://technet.microsoft.com/en-us/library/hh831823\(v=ws.11\).aspx](https://technet.microsoft.com/en-us/library/hh831823(v=ws.11).aspx). Accessed: 2023-01.
- [25] Zhong-Zhen Wu and Han-Chiang Chen. Design and implementation of tcp/ip offload engine system over gigabit ethernet. In *Proceedings of 15th International Conference on Computer Communications and Networks*, pages 245–250. IEEE, 2006.
- [26] Krishna Kant. Tcp offload performance for front-end servers. In *GLOBECOM’03. IEEE Global Telecommunications Conference (IEEE Cat. No. 03CH37489)*, volume 6, pages 3242–3247. IEEE, 2003.
- [27] Douglas Freimuth, Elbert C Hu, Jason D LaVoie, Ronald Mraz, Erich M Nahum, Prashant Pradhan, and John M Tracey. Server network scalability and tcp offload. In *USENIX Annual Technical Conference, General Track*, pages 209–222, 2005.
- [28] Ruining Chen and Guoao Sun. A survey of kernel-bypass techniques in network stack. In *Proceedings of the 2018 2nd International Conference on Computer Science and Artificial Intelligence*, pages 474–477, 2018.
- [29] 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. *ACM SIGCOMM Computer Communication Review*, 45(4):537–550, 2015.

-
- [30] Gautam Kumar, Nandita Dukkipati, Keon Jang, Hassan MG Wassel, Xian Wu, Behnam Montazeri, Yaogong Wang, Kevin Springborn, Christopher Alfeld, Michael Ryan, et al. Swift: Delay is simple and effective for congestion control in the datacenter. In *Proceedings of the Annual conference of the ACM Special Interest Group on Data Communication on the applications, technologies, architectures, and protocols for computer communication*, pages 514–528, 2020.
- [31] Manage virtual machines with virt-manager. <https://virt-manager.org>, . Accessed: 2023-03.
- [32] Oracle vm virtualbox. <https://www.virtualbox.org/>, . Accessed: 2023-03.