

# Guest-Based System Call Introspection with Extended Berkeley Packet Filter

by

*Huzaiifa Patel*

A thesis proposal submitted to the School of Computer Science in partial fulfillment  
of the requirements for the degree of

**Bachelor of Computer Science**

Under the supervision of Dr. Anil Somayaji

Carleton University

Ottawa, Ontario

September, 2022

© 2022 Huzaiifa Patel

*their kindness is masquerade.*

*yearning to occupy one with false pretenses.*

*it's used to sedate.*

*I promise you'll get this when the sky clears for you.*

# Abstract

## Acknowledgments

I want to express my heartfelt gratitude to my supervisor, Dr. Anil Somayaji for providing me with the opportunity to work on a thesis during the final year of my undergraduate degree. Unlike previous variations of the Computer Science undergraduate degree requirements, completing a thesis is no longer a prerequisite. Therefore, I prostualte it is a great privlidge and honor to be given the opportunity to enroll into a thesis-based course during ones undergraduate studies.

I did not have prior experience in formal research when I first approached Dr. Somayaji. Despite this shortcoming, it did not stop him from investing his time and resources towards my academic growth. Without his feedback and ideas on my framework implementation and writing of this thesis, as well as his expertise in eBPF, Hypervisors, and Unix based Operating Systems, this thesis would not have been possible.

I would like to commend PhD student Manuel Andreas from The Technical University of Munich for introducing me to the concept of a Hypervisor. Without him, I would not have approached Dr. Somayaji with the intention of wanting to conduct research on them. His minor action of introducing me to hypervisors had the significant effect of inspiring me to write a thesis on the subject. I also want to thank him for his willingness to endlessly and tirelessly teach, discuss and help me understand the intricacies of hypervisors, the Linux kernel, and the C programming language.

I would also like to thank Carleton University's faculty of Computer Science for their efforts in imparting knowledge that has enthraled and inspired me to learn all that I can about Computer Science.

I would like to extend my appreciation to the various internet communities which have provided the world with invaluable compiled resources on hypervisors, Unix based operating systems, eBPF, the Linux kernel, the C programming language, and Latex, which has helped me tremendously in writing this thesis.

Finally, I would like to thank my immediate family for their encouragement and support towards my research interests and educational pursuits.

# Contents

<b>Abstract</b>	<b>i</b>
<b>Acknowledgments</b>	<b>ii</b>
<b>Nomenclature</b>	<b>vi</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Motivation . . . . .	2
1.1.1 Why Design a New Framework? . . . . .	4
1.1.2 Why Out-Of-VM monitor? . . . . .	4
1.1.3 Why eBPF? . . . . .	4
1.1.4 Why Sequences of System Calls? . . . . .	4
1.2 Problem . . . . .	4
1.2.1 The Semantic Gap Problem . . . . .	4
1.3 Approaching the Problem . . . . .	4
1.4 Contributions . . . . .	4
1.5 Thesis Organization . . . . .	4
<b>2 Background</b>	<b>5</b>
2.1 Hypervisor . . . . .	5
2.2 Intel Virtualization Extention (VT-X) . . . . .	6
2.3 The Kernel Virtual Machine Hypervisor . . . . .	6
2.3.1 Model Specific Registers . . . . .	7
2.3.2 VMCS . . . . .	7
2.3.3 VM ENTRY Context Switch . . . . .	7
2.3.4 VM EXIT Context Switch . . . . .	7
2.4 QEMU . . . . .	7

2.5	System Calls . . . . .	7
2.6	Virtual Machine Introspection . . . . .	7
2.7	eBPF . . . . .	7
2.8	The Linux Kernel Tracepoint API . . . . .	7
2.9	pH-based Sequences of System Call . . . . .	7
<b>3</b>	<b>Related work</b>	<b>8</b>
3.1	Nitro: Hardware-Based System Call Tracing for Virtual Machines . . . . .	8
<b>4</b>	<b>Designing Frail</b>	<b>9</b>
<b>5</b>	<b>Implementing Frail</b>	<b>10</b>
5.1	User Space Component . . . . .	10
5.2	Kernel Space Component . . . . .	10
5.2.1	Custom Linux Kernel Tracepoint . . . . .	10
5.2.2	Kernel Module . . . . .	10
5.3	Tracing Processess . . . . .	10
5.4	Proof of Tracability of all KVM Guest System Calls . . . . .	10
<b>6</b>	<b>Threat Model of Frail</b>	<b>11</b>
<b>7</b>	<b>Future Work</b>	<b>12</b>
<b>8</b>	<b>Conclusion</b>	<b>13</b>
<b>9</b>	<b>References</b>	<b>14</b>

## Nomenclature

<b>VM</b>	Virtual Machine
<b>KVM</b>	Kernel-based Virtual Machine
<b>OS</b>	Operating System
<b>VMI</b>	Virtual Machine Introspection
<b>CPU</b>	Central Processing Unit
<b>AMD-V</b>	Advanced Micro Devices Virtualization
<b>VT-x</b>	Intel Virtualization Extension
<b>MSR</b>	Model Specific Register
<b>VMM</b>	Virtual Machine Monitor
<b>EFER</b>	Extended Feature Enable Register
<b>eBPF</b>	Extended Berkeley Packet Filter



# Introduction

Cloud computing has been ever-increasing in demand due to its scalability, security and convenience characteristics. These trends can largely be attributed to the usage of virtualization. Fundamentally, virtualization is a technology that makes it possible for multiple different operating systems (OSs) to run concurrently, and in an isolated environment on the same hardware. It makes use of a machine's hardware to support the software that creates and manages virtual machines (VMs). The software that creates and manages virtual machines has two formal names: (1) hypervisor and (2) virtual machine monitor (VMM). The operating system running a hypervisor is called the host OS, while the virtual machine that uses its resources is known as the guest OS.

The fundamental reason for introducing a hypervisor layer on a machine is that before hypervisors, a machine could run only one operating system at a time. This constraint often led to wasted resources, as a single OS seldom fully utilized the hardware's capacity. The computing capacity of a CPU is so large that it is difficult for one operating system to efficiently use all of it. Hypervisors address the above constraint by aggregating the resources of virtualized physical servers (such as memory, network bandwidth and CPU cycles) and then allocating those resources to virtual machines. This allowed users to switch between one machine, many operating systems, and multiple applications at their discretion.

Since hardware-assisted virtualization was introduced to commodity x86 servers ten years ago, it has become the common practice for server deployment. Today, about 75% of x86 server workloads run in virtual machines (VMs). Recent researchers classified the virtualization market to be composed of four main participants (taking up to 93% of the total market share) which are: VMware, Microsoft's Hyper-V, Xen and KVM.

These first two brands are the main two open-source in the market. VMware is also the most known, offering a total presence of 81% and 52% as primary, Xen is second with 81% and 18% as primary, KVM followed by 58% presence and 9% primary and Hyper-V by 43% presence and 9% primary.

As such, the role of a hypervisor is highly security critical. The successful exploitation of a hypervisor can result in a complete breach of isolation between the clients, resulting in loss of availability of client services, non-public information becoming accessible to unauthorized parties, and data, software or hardware being altered by unauthorized parties. Because of the above possible effects of hypervisor exploitation, effective methodologies for monitoring hypervisors are required.

This thesis is a first step towards fixing the shortcomings of not being able to detect and respond to malicious anomalies. We present a framework named Frail, a KVM hypervisor and Intel VT-x exclusive out-of-vm system call monitoring, detection and response system. Our framework is implemented using a combination of already existing software. Firstly, it uses Extended Berkeley Packet Filter (eBPF) to extract system call and guest OS process information. Secondly, it uses pH's implementation of sequences of system calls to detect. Lastly, it uses our own implementation to respond to malicious anomalies by slowing down the guest process responsible for the detected anomaly. In this thesis, we will examine the design and implementation of Frail, explore its security implication on the KVM hypervisor and its guests, and explore its impact on system performance.

## 1.1 Motivation

<https://dl.acm.org/doi/pdf/10.1145/2815400.2815420>: Since hardware-assisted virtualization was introduced to commodity x86 servers ten years ago, it has become the common practice for server deployment [7]. Today, about 75server workloads run in virtual machines (VMs) [13]. Virtualization enables the consolidation of multiple VMs on a single server, thereby reducing hardware and operation costs [14]. Virtualization promises to reduce these costs without sacrificing robustness and security. We contend, however, that this promise is not fulfilled in practice, because hypervisors—the software

layers that run VMs—are bug-prone. Hypervisor bugs can cause an operating system (OS) that runs within a VM to act incorrectly, crash, or become vulnerable to security exploits [18]. Hypervisor bugs are software bugs, but the damage they cause is similar to that of hardware bugs. Since hypervisors virtualize the hardware of VMs, their bugs cause the VMs to experience that the underlying hardware violates its specification. Patching hypervisor bugs is much easier than fixing the hardware, yet doing so may induce VM downtime and deter cloud customers, as indeed experienced by leading cloud providers [24, 71].

Current computer systems have no general-purpose mechanism for detecting and responding to successful exploitation of the KVM hypervisor. We can't rely on users to detect and respond to exploits of the KVM hypervisor because as our computer systems continue to grow increasingly complex, so too does it become more difficult to measure precisely what they are doing at any given moment. As a result, users often have a limited notion of what is happening on their computer systems internal states. An unfortunate consequence of the lack of selfawareness in this domain is that it decreases the likelihood of spotting and appropriately reacting to malicious anomalies. If users are unable to adequately monitor our computer systems, computers should be programmed to watch over themselves.

- 1.1.1 Why Design a New Framework?
- 1.1.2 Why Out-Of-VM monitor?
- 1.1.3 Why eBPF?
- 1.1.4 Why Sequences of System Calls?
- 1.2 Problem
- 1.2.1 The Semantic Gap Problem
- 1.3 Approaching the Problem
- 1.4 Contributions
- 1.5 Thesis Organization

# Background

## 2.1 Hypervisor

There are two ways a hypervisor can virtualize a machine:

A hypervisor runs Guest OS instructions either directly on the host's CPU, or on the host OS. In both scenarios, the goal of a hypervisor is to provide a software-controlled layer that resembles the host hardware. Hypervisors can be classified into two types that are dependent on how they to runs Guest OS instructions.

(1) Type 1 (bare metal) hypervisors, which runs Guest OS instructions directly on the host's hardware in order to control the hardware and monitor the guest OS. Typical examples of such hypervisors include Xen, VMware ESX, and Microsoft Hyper-V.

(2) Type 2 (hosted) hypervisors, which run within a traditional OS. In other words, a hosted hypervisor adds a distinct software layer atop the host OS, and the guest OS becomes a third software layer above the hardware. Well-known examples of type 2 hypervisors include KVM, VMware Workstation, VirtualBox, and QEMU.

Although the preceding type 1 and type 2 hypervisor classification has been widely accepted, it is not clear it insufficiently differentiates among hypervisors of the same type (e.g., KVM vs. QEMU).

KVM is not a clear case as it could be categorized as either one. The KVM kernel module turns Linux kernel into a type 1 bare-metal hypervisor, while the overall system could be categorized to type 2 because the host OS is still fully functional and the other VM's are standard Linux processes from its perspective.

There-fore, based on how the virtualization gets designed (hardware vs. software) and the guest OS and its application code is executed, we can have another type of classification of hypervisors that will be used throughout this thesis:

(1) Native hypervisors that directly push the guest code to execute natively on the hardware using hardware virtualization.

(2) Emulation hypervisors that translate each guest instruction for an emulated execution using software virtualization.

Examples of native hypervisors include Xen, KVM, VMware ESX, and Microsoft HyperV, and emulation hypervisors include QEMU, Bochs, and the very early versions of VMware-Workstation and VirtualBox (note that recent VMware-Workstation and VirtualBox are able to execute the guest OS code natively). Since there is no binary code translation involved, native hypervisor runs much faster than emulation hypervisor.

In this thesis, we will be solely on the KVM VM.

Hardware-assisted

## **2.2 Intel Virtualization Extention (VT-X)**

## **2.3 The Kernel Virtual Machine Hypervisor**

Kernel-based Virtual Machine (KVM) is a hypervisor that is implemented as a Linux kernel module that allows the kernel to function as a hypervisor. It was merged into the mainline Linux kernel in version 2.6.20, which was released on February 5, 2007. KVM requires a CPU with hardware virtualization extensions, such as Intel VT-x or AMD-V. While working with KVM, we will only be focusing on Intel VT-x hardware virtualization.

**2.3.1 Model Specific Registers**

**2.3.2 VMCS**

**2.3.3 VM ENTRY Context Switch**

**2.3.4 VM EXIT Context Switch**

**2.4 QEMU**

**2.5 System Calls**

**2.6 Virtual Machine Introspection**

**2.7 eBPF**

**2.8 The Linux Kernel Tracepoint API**

**2.9 pH-based Sequences of System Call**

## Related work

There are 12 projects that use the guest-assisted approach. The pioneer work, LARES [Payne et al. 2008], inserts hooks in a guest VM and protects its guest component by using the hypervisor for memory isolation with the goal of supporting active monitoring. Unlike passive monitoring, active monitoring requires the interposition of kernel events. As a result, it requires the monitoring code to be executed inside the guest OS, which is why it essentially leads to the solution of inserting certain hooks inside the guest VM. The hooks are used to trigger events that can notify the hypervisor or redirect execution to an external VM. More specifically, LARES design involves three components: a guest component, a secure VM, and a hypervisor. The hypervisor helps to protect the guest VM component by memory isolation and acts as the communication component between the guest VM and the secure VM. The secure VM is used to analyze the events and take actions necessary to prevent attacks.

### 3.1 Nitro: Hardware-Based System Call Tracing for Virtual Machines



## Designing Frail

## **Implementing Frail**

### **5.1 User Space Component**

### **5.2 Kernel Space Component**

#### **5.2.1 Custom Linux Kernel Tracepoint**

#### **5.2.2 Kernel Module**

### **5.3 Tracing Processess**

### **5.4 Proof of Tracability of all KVM Guest System Calls**

## Threat Model of Frail

## Future Work

## Conclusion

## References