



**Hochschule für Technik
und Wirtschaft Berlin**

University of Applied Sciences

System Calls for Containerising and Managing Processes in Linux

A Thesis

Submitted in Partial Fulfillment of the Requirements for the Degree of

Master of Science (M.Sc.) in Applied Computer Science

at the

Berlin University of Applied Sciences (HTW)

First Supervisor: Prof. Dr. Sebastian Bauer

Second Supervisor: Dr. Michael Witt

Author: B.Sc. Atanas Denkov

Matriculation Number: s0559025

Submission Date: 12.07.2022

Contents

1	Introduction	1
1.1	Motivation	1
1.2	Objectives	1
1.3	Content Structure	1
2	Fundamentals	2
2.1	Virtualisation	2
2.1.1	Axioms	2
2.1.2	Hardware Virtualisation	3
2.1.3	Operating System Virtualisation	4
	List of Figures	I
	List of Tables	II
	Source Code Content	III
	References	IV

Chapter 1

Introduction

1.1 Motivation

Primitive support for multiprocessing in the form of basic context switching and dedicated input-output components was introduced in the late 1950s. Multiprocessing allowed for concurrent execution of multiple instructions at the cost of increased system complexity. Interleaved processes had a global unrestricted view of the system which inevitably led to unpredictable program behaviour. For example, programs had the ability to modify each other's memory and monopolise computer resources. Hence, to ensure correctness, every program had to carefully manage its interactions with hardware and all other processes in the system, which resulted in an unsustainable programming model.

The aforementioned issues were addressed by shifting the responsibility of resource management and process protection into a privileged control program that acted as an intermediary between hardware and user programs.

1.2 Objectives

1.3 Content Structure

Chapter 2

Fundamentals

2.1 Virtualisation

2.1.1 Axioms

Noninterference

Codd et al. [Cod59] summarise the fundamental requirements of a multiprogramming system and emphasise the concept of noninterference between processes across space and time. *Spatial noninterference* is represented by all mechanisms that protect references to memory, disk and input-output devices [Cod59]. For example, memory segmentation is a method found in operating system kernels that assigns each process a dedicated portion of physical memory that is invisible to all other processes in the system. The kernel traps any attempt made by a process to access memory outside its allocated memory segment, thereby guaranteeing spatial noninterference [SGG18]. *Temporal noninterference* refers to those mechanisms that allocate execution time and protect against the monopolisation thereof [Cod59]. For instance, CPU scheduling is a technique that decides which process shall run on a core such that the core does not idle and all processes make sufficient runtime progress [SGG18]. The scheduling semantics, paired with an interrupt mechanism that makes sure that no process has hold of the core for too long, guarantee temporal noninterference.

Isolation

Anjali, Caraza-Harter, and Swift [ACS20] define isolation as the level of dependency that a virtualisation platform has towards the host kernel. We generalise this definition and say that *isolation* is the level of dependency that one piece of software has to another. Conceptually, isolation deals with explicit vertical relationships between software, and noninterference deals with implicit horizontal relationships between processes. The level of dependency can be quantified by counting the lines of external code that a software executes to obtain a desired functionality.

For example, Anjali, Caraza-Harter, and Swift [ACS20] count the lines of kernel code that a virtualisation platform executes when providing services to sandboxed applications. High counts indicate a strong dependency, i.e weak isolation.

Performance

Randal [Ran20] defines performance as the contention between the overhead associated with isolating a process from its environment and the benefits of sharing resources between processes, i.e fully utilising the capacity of the underlying resource pool. Anjali, Caraza-Harter, and Swift [ACS20] use a similar definition and measure the level of isolation provided by three different virtualisation platforms contrasted against processing unit, memory and input-output performance metrics. Isolation is measured as described in Chapter 2.1.1. The lines of code are mapped to the kernel subsystem in which they reside, resulting in an isolation metric per subsystem. Workloads are then executed as sandboxed applications, instantiated multiple times to measure resource utilisation on a single host system, and compared against the lines of code executed by the subsystem that the workload corresponds to. For example, the authors define an application that computes prime numbers up to a limit. Since the workload is compute-bound, processing speed is measured and compared to the number of lines of executed code that reside in the processing unit and architecture-specific subsystems of the kernel.

Manco et al. [Man17] use *same-host density* as a performance metric that measures the number of sandboxed applications that can be consolidated onto a single server.. In addition, *boot*, *pause* and *unpause times* are also considered as important performance indicators for particular use cases, such as elastic content delivery networks [Kue17] [Man17]. The authors “sequentially start 1000 virtual machines and measure the time it takes to create each VM and the time it takes the VM to boot” [Man17, p. 5].

2.1.2 Hardware Virtualisation

Popek and Goldberg [PG74] refer to the control program as a *virtual machine monitor* that ensures isolation and noninterference by providing every program with an environment that is “[...] effect identical with that demonstrated if the program had been run on the original machine directly” [PG74, p. 2]. This definition implies that a running program does not directly use the bare metal resources available. Instead, resources are emulated by the virtual machine monitor at the instruction level and presented as a dedicated physical system. Such an environment is called a *virtual machine*.

Popek and Goldberg [PG74] define a requirement that the instruction-set architecture of a computer has to satisfy for it to be virtualisable. The instruction set must be segregated into three groups of instructions - privileged, sensitive and innocuous. An instruction is privileged if it

requires changing the mode of execution from user to supervisor mode by means of a trap [PG74]. An instruction i is control-sensitive if, when applied to the current processor state S_1 , results in a new state $i(S_1) = S_2$ such that the execution mode of S_2 does not equal that of S_1 or if S_2 has access to different resources than S_1 or both [PG74]. An instruction is behaviour-sensitive if its execution depends on the execution mode or its position in memory [PG74]. An instruction is innocuous if it is not sensitive. Given these definitions, a computer is virtualisable “[...] if the set of sensitive instructions for that computer is a subset of the set of privileged instructions” [PG74, p. 6]. If this criterion is met, the virtual machine monitor can trap all sensitive instructions and emulate each via a homomorphism $i : C_r \rightarrow C_v$ that maps the state space of the processor without the virtual machine monitor loaded C_r to the state space with the virtual machine monitor loaded C_v [PG74]. Innocuous instructions do not require protection, i.e a homomorphic mapping, and are directly executed by the processor [PG74].

The aforementioned homomorphism enables a virtual machine to host a kernel, referred to as a *guest kernel*, that “believes” it is interacting with a physical system. Consequently, the isolation and noninterference boundaries between user programs running on different virtual machines are stronger compared to user processes running on a shared kernel. Even if a guest kernel becomes compromised or encounters an unrecoverable error condition, other virtual machines remain unaffected.

The performance cost of hardware virtualisation becomes apparent when measuring same-host density and boot times. Manco et al. [Man17] consider memory consumption and on-disk image size as the primary limiting factors. Lv et al. [Lv12] come to the same conclusion and further show that the high frequency of context switches also has a negative impact.

2.1.3 Operating System Virtualisation

List of Figures

List of Tables

Source Code Content

References

- [ACS20] Anjali, Tyler Caraza-Harter, and Michael M. Swift. “Blending Containers and Virtual Machines: A Study of Firecracker and GVisor”. In: *Proceedings of the 16th ACM SIGPLAN/SIGOPS International Conference on Virtual Execution Environments*. VEE ’20. Lausanne, Switzerland: Association for Computing Machinery, 2020, pp. 101–113. ISBN: 9781450375542. DOI: 10.1145/3381052.3381315. URL: <https://doi.org/10.1145/3381052.3381315>.
- [Cod59] E. F. Codd et al. “Multiprogramming STRETCH: Feasibility Considerations”. In: *Commun. ACM* 2.11 (Nov. 1959), pp. 13–17. ISSN: 0001-0782. DOI: 10.1145/368481.368502. URL: <https://doi.org/10.1145/368481.368502>.
- [Kue17] Simon Kuenzer et al. “Unikernels Everywhere: The Case for Elastic CDNs”. In: *Proceedings of the 13th ACM SIGPLAN/SIGOPS International Conference on Virtual Execution Environments*. VEE ’17. Xi’an, China: Association for Computing Machinery, 2017, pp. 15–29. ISBN: 9781450349482. DOI: 10.1145/3050748.3050757. URL: <https://doi.org/10.1145/3050748.3050757>.
- [Lv12] Hui Lv et al. “Virtualization Challenges: A View from Server Consolidation Perspective”. In: *Proceedings of the 8th ACM SIGPLAN/SIGOPS Conference on Virtual Execution Environments*. VEE ’12. London, England, UK: Association for Computing Machinery, 2012, pp. 15–26. ISBN: 9781450311762. DOI: 10.1145/2151024.2151030. URL: <https://doi.org/10.1145/2151024.2151030>.
- [Man17] Filipe Manco et al. “My VM is Lighter (and Safer) than Your Container”. In: *Proceedings of the 26th Symposium on Operating Systems Principles*. SOSP ’17. Shanghai, China: Association for Computing Machinery, 2017, pp. 218–233. ISBN: 9781450350853. DOI: 10.1145/3132747.3132763. URL: <https://doi.org/10.1145/3132747.3132763>.
- [PG74] Gerald J. Popek and Robert P. Goldberg. “Formal Requirements for Virtualizable Third Generation Architectures”. In: *Commun. ACM* 17.7 (July 1974), pp. 412–421. ISSN: 0001-0782. DOI: 10.1145/361011.361073. URL: <https://doi.org/10.1145/361011.361073>.
- [Ran20] Allison Randal. “The Ideal Versus the Real: Revisiting the History of Virtual Machines and Containers”. In: *ACM Comput. Surv.* 53.1 (Feb. 2020). ISSN: 0360-0300. DOI: 10.1145/3365199. URL: <https://doi.org/10.1145/3365199>.

-
- [SGG18] Abraham Silberschatz, Peter B. Galvin, and Greg Gagne. *Operating System Concepts*. 10th. Wiley Publishing, 2018. ISBN: 9781119320913.