

University of Warsaw
Faculty of Mathematics, Informatics and Mechanics

Krzysztof Małysa

Student no. 394442

Multi-process sandbox for unprivileged users on Linux

Master's thesis
in COMPUTER SCIENCE

Supervisor:
dr Janina Mincer-Daszkiewicz
Institute of Informatics

Warsaw, December 2022

Abstract

TODO

Keywords

sandboxing, security, Linux, secure execution, arbitrary code execution

Thesis domain (Socrates-Erasmus subject area codes)

11.3 Informatics, Computer Science

Subject classification

Security and privacy – Systems security – Operating systems security

Tytuł pracy w języku polskim

Sandbox wielu procesów dla nieuprzywilejowanych użytkowników systemu Linux

Contents

1. Introduction	5
1.1. Background	5
1.2. Goal of the thesis	5
1.2.1. Requirements	6
1.2.2. Existing solutions	6
1.3. Scope	7
1.4. Significance	7
1.5. Structure of the Thesis	7
2. Literature overview	9
3. Methodology	11
4. Design	13
5. Implementation	15
6. Performance Evaluation	17
7. Use Cases and Applications	19
8. Integration with Online Judge	21
9. Future work and Opportunities	23
10. Conclusion	25

Chapter 1

Introduction

1.1. Background

Secure execution environments are commonplace these days, from containers and virtual machines on servers to sandboxes on laptop and smartphones — most of which run on Linux. They are used to securely execute untrusted code, as well as trusted programs to prevent damage escalation in the event of unknown vulnerabilities. Their key features are isolation, limiting resource usage, and accounting for resource consumption.

The features of Linux allow the creation of simple yet effective and efficient secure environments. They work at application runtime, so in most cases existing software does not need to be adapted to use them. This makes them easily applicable, and explains why they adoption is growing.

In this thesis, the most important application of sandboxing are online judge systems. Online judge systems have beneficial role in programming education and competitive programming. They allow testing user-provided solution to a specific problem. The solution is run on a predefined test cases in order to check if it is valid. In such platforms isolating the compilation and running of the solution is essential to provide security and robustness of the platform itself.

Historically, isolation techniques evolved together with the online judge platforms. The most primitive (yet insecure) was usage of `chroot(2)` [10] to restrict access to part of the filesystem. To increase isolation virtual machines were used [2]. Later, containerization became a new way to provide isolation [5, 13].

Online education platforms greatly facilitate teaching and learning programming. They provide quick feedback on the correctness of the code the user submits. They are used in schools and universities and provide great learning opportunities for all.

1.2. Goal of the thesis

The goal of the thesis is to design, implement and integrate a new sandbox for the Sim project [3]. The Sim project [3] is an online platform for preparing people for and carrying out algorithmic contests. It has an online judge with a sandbox specially developed for this use case. Over the years the sandbox became a limitation. It only allows running a single-threaded statically linked executable of programs written in C, C++ or Pascal. The new sandbox will allow supporting more programming languages and improve security of the solution compilation stage.

1.2.1. Requirements

The new sandbox needs to be optimized for running short-running programs as well as have minimal runtime overhead. Most of the test cases the solution is run on are small and solution completes them in less than 10ms. The goal is to allow hundreds of such sort-running runs per second, hence optimizing for short-running programs is important. However, minimizing overhead of the sandbox during the run is also important i.e. if the program runs X ms normally, the objective is that the program inside the new sandbox will also run approximately X ms.

The new sandbox needs to be versatile. It will be used to secure the compilation of the solutions as well as running of the solutions. Compilation is a complicated process that involves parsing, translating, optimizing and linking the final program. It often requires running several executables e.g. compiler and linker, so allowing a single process inside sandbox is not enough. Sandboxing solution is simpler, because it is a single process. But since it is often short-running, the overhead needs to be minimal.

The sandbox needs to allow limiting resources. Real time, CPU time, memory – these need to be limited not only for the robustness of the platform, but specific problems require different limits. The goal of some problems is to solve it with very restricted memory e.g. find a missing integer in a random permutation of integers $1, \dots, n$ without one element, but in constant memory.

The sandbox needs to account resource usage. For every test, the user is presented with consumed memory and CPU time by their solution. The sandbox needs to provide this information.

The last requirement is the sandbox will not require any privileges. There is a tool called Sip [4] for preparing the problem packages for the Sim platform. One of the purposes of the tool is to run the solutions inside the same secure environment as on the Sim platform. The user should not need any privileges to run this tool, so the sandbox should not require them either.

1.2.2. Existing solutions

Approaches to form a secure execution environments differ. One of them is virtualization or emulation e.g. QEMU [8] and KVM [7], VirtualBox [9], VMWare Workstation [15]. Although powerful and effective, they come with an enormous overhead i.e. booting up an entire operating system. Moreover, emulation noticeably slows down the runtime of an emulated application, rendering such solutions inapplicable.

Containers provide much lower overhead: setup of an order of milliseconds and negligible runtime overhead. But, Docker [6], LXC [1] require root privileges to create a container. systemd-nspawn [14] requires root privileges to run.

Rootless are containers [11] that can be created and run by an unprivileged user are the almost perfect solution to the problem. They provide almost all of the functionality of the normal containers but without the need to engage a privileged user. However, they often use `setuid` binaries and that is undesirable [12]. Also they are not optimized to run sequences of short-running programs. In this thesis we will create a sandbox that uses the same techniques as rootless containers but will be optimized for running sequences of short-running programs.

1.3. Scope

The sandbox is implemented in C++ and runs only on Linux. The implementation was tested on x86_64 but should also work on x86 or arm64, due to no architecture specific parts. Only seccomp filters may need to be adapted for other architectures, but the sandbox itself uses user-provided seccomp filters so it should work out of the box.

The sandbox uses linux namespaces to isolate the sandboxed program and cgroups to limit resources and provide resource usage accounting.

A sandbox client in Rust language would be beneficial for the broader adoption of the sandbox, but is intentionally out of scope of this thesis.

1.4. Significance

A highly efficient sandbox could improve user experience by reducing time the user waits for the program to be judged by the online judge platform and could enable supporting more programming languages.

The versatile sandbox can also be used to e.g. sanitize compiling a PDF from L^AT_EX sources. It could also be used for safely executing untrusted server-side scripts in web applications.

1.5. Structure of the Thesis

Chapter 2 contains overview of sandboxing techniques and existing implementations and comparative analysis of them. Chapter 3 describes briefly the all aspects of the development process from architecture and design to testing and debugging and performance evaluation. Details of design and architecture are described in chapter 4. Implementation is described in 5. Chapter 6 contains performance evaluation of the final implementation and impact of some optimizations. Later, in chapter 7 use cases and applications are discussed. Chapter 8 details integration with online judge and challenges involved. In chapter 9 future work and opportunities are discussed. Finally, chapter 10 contains the conclusion.

Chapter 2

Literature overview

Chapter 3

Methodology

Chapter 4

Design

Chapter 5

Implementation

Chapter 6

Performance Evaluation

Chapter 7

Use Cases and Applications

Chapter 8

Integration with Online Judge

Chapter 9

Future work and Opportunities

Chapter 10

Conclusion

Bibliography

- [1] David Beserra et al. “Performance Analysis of LXC for HPC Environments.” In: *CISIS*. IEEE Computer Society, 2015, pp. 358–363. ISBN: 978-1-4799-8870-9. URL: <http://dblp.uni-trier.de/db/conf/cisis/cisis2015.html#BeserraMEBSF15>.
- [2] Sander van der Burg and Eelco Dolstra. “Automating System Tests Using Declarative Virtual Machines”. In: *2010 IEEE 21st International Symposium on Software Reliability Engineering*. 2010, pp. 181–190. DOI: [10.1109/ISSRE.2010.34](https://doi.org/10.1109/ISSRE.2010.34).
- [3] Krzysztof Małysa. *Sim project*. URL: <https://github.com/varqox/sim> (visited on 03/15/2023).
- [4] Krzysztof Małysa. *Sip – a tool for preparing problem packages for the Sim platform*. URL: <https://github.com/varqox/sip> (visited on 03/15/2023).
- [5] Martin Mareš and Bernard Blackham. “A New Contest Sandbox.” In: *Olympiads in Informatics* 6 (2012).
- [6] Dirk Merkel. “Docker: Lightweight Linux Containers for Consistent Development and Deployment”. In: *Linux J*. 2014.239 (Mar. 2014). ISSN: 1075-3583. URL: <http://dl.acm.org/citation.cfm?id=2600239.2600241>.
- [7] *Official website of Kernel Virtual Machine*. URL: <https://www.linux-kvm.org/> (visited on 11/23/2022).
- [8] *Official website of QEMU — A generic and open source machine emulator and virtualizer*. URL: <https://www.qemu.org/> (visited on 11/23/2022).
- [9] Oracle. *Official website of VirtualBox*. URL: <https://www.virtualbox.org/> (visited on 11/23/2022).
- [10] Vassilis Prevelakis and Diomidis Spinellis. “Sandboxing Applications.” In: *Usenix annual technical conference, freenix track*. Citeseer. 2001, pp. 119–126.
- [11] rootlesscontainers.rs. *Rootless Containers*. URL: <https://rootlesscontainers.rs> (visited on 11/28/2022).
- [12] Giuseppe Scrivano. *Rootless containers with Podman and fuse-overlayfs*. June 4, 2019. URL: https://indico.cern.ch/event/757415/contributions/3421994/attachments/1855302/3047064/Podman_Rootless_Containers.pdf (visited on 11/28/2022).
- [13] František Špaček, Radomír Sohlich, and Tomáš Dulík. “Docker as Platform for Assignments Evaluation”. In: *Procedia Engineering* 100 (2015). 25th DAAAM International Symposium on Intelligent Manufacturing and Automation, 2014, pp. 1665–1671. ISSN: 1877-7058. DOI: <https://doi.org/10.1016/j.proeng.2015.01.541>. URL: <https://www.sciencedirect.com/science/article/pii/S1877705815005688>.

- [14] systemd. *systemd-nspawn* — *Spawn a command or OS in a light-weight container*. URL: <https://www.freedesktop.org/software/systemd/man/systemd-nspawn.html> (visited on 11/28/2022).
- [15] VMWare. *Official website of VMWare Workstation*. URL: <https://www.vmware.com/products/workstation/> (visited on 11/23/2022).