# Sandbox for multi-process applications for unprivileged users on Linux

Krzysztof Małysa

University of Warsaw

December 2023

# The Sim platform

# Before the new sandbox

## Old sandbox

- Single-threaded, statically-linked executables
- C, C++, Pascal
- Overhead — `ptrace`

## Problems

- No support for other languages e.g. Python
- Compilation

## Solution

A new sandbox.

# The submission on the Sim platform

## Requirements for new Sandbox

1. Versatile
2. Support for multi-process applications
3. Low overhead
4. Optimized for short-running programs
5. Limiting resources
6. Runtime statistics
7. For unprivileged users

- Often require privileges e.g. OS modifications.
- Some use `ptrace`. Slow, TOCTOU problem.
- Few provide runtime statistics.
- None is optimized for short-running programs.

# Design: client-server

Allows sharing resources between requests.



Benchmarks in a moment.

# Used Linux kernel mechanisms

Linux namespaces, cgroups, `prlimit`, seccomp BPF filters.

# Benchmark of some performance optimizations

| Benchmark | Mean request time | Std. dev. | Std. err. on the mean | Slowdown |
|---|---|---|---|---|
| Baseline | 2.348ms | 0.768ms (32.71%) | 0.024ms (1.03%) | 0.00% |
| New network namespace for each request | 2.970ms | 0.856ms (28.83%) | 0.027ms (0.91%) | 26.49% |
| New IPC namespace for each request | 2.522ms | 0.782ms (31.02%) | 0.025ms (0.98%) | 7.41% |
| New UTS namespace for each request | 2.478ms | 0.771ms (31.14%) | 0.024ms (0.98%) | 5.54% |

Table: Statistics for each row were collected from 1000 runs. Each row contains real time it took to handle request to sandbox the /bin/true program.

| Sandbox | Mean time | Std. dev. | Std. err. on the mean | Slowdown |
|---|---|---|---|---|
| no sandbox | 0.893ms | 0.409ms (45.80%) | 0.013ms (1.45%) | 1x |
| sandbox | 2.348ms | 0.768ms (32.71%) | 0.024ms (1.03%) | 2.39x |
| nsjail | 10.393ms | 1.327ms (12.77%) | 0.042ms (0.40%) | 10.57x |

Table: Statistics for each row were collected from 1000 runs. Each row contains real time it took to handle request to sandbox the /bin/true program. While the slowdown of the sandbox is huge (more than twofold), it still allows for hundreds of runs per second and that was the goal of this thesis, whereas nsjail is more than 4 times slower than our sandbox.

Tested compilation of all solutions of 4 problems from finals of the XXII Polish Olimpiad in Informatics.

- Compilation.
- Running the model solution on the tests.

### Benchmark

Collect statistic from 10 runs.

Tested compilation of all solutions of 4 problems from finals of the XXII Polish Olimpiad in Informatics. $18 + 18 + 13 + 8 = 57$ solutions.

- Only some compilations present statistically significant difference.
- Half of the time it is statistically no slower that without the sandbox.
- 0% – 24% slowdown. Most of the time below 10%.
- Almost always seccomp BPF filters cause the overhead (up to 15% without it).
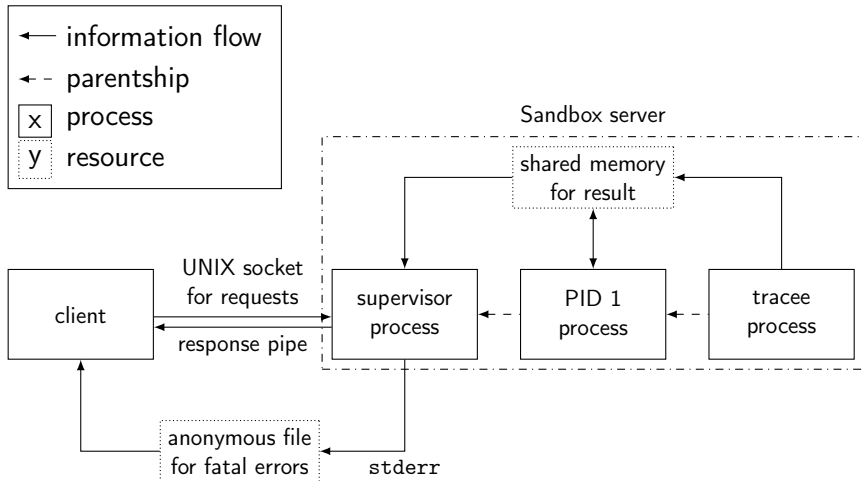
For each problem the single model solution was tested. Problems have: $55 + 57 + 69 + 65 = 246$ tests.

- Vast majority of tests shows no statistical difference.
- Few test show slowdown — up to 19%.
- Medium tests even show speed-up — up to 44%.
- Low number of system calls. Most of the time, computations without IO.

# Future work

- Setting CPU affinity
- Support for networking
- Rust frontend / client
- Further experimentation

Thank you.

1. Time limits
2. Runtime statistics
3. Error handling (stderr)
4. Request sending and receiving (serialization)
5. File descriptors
6. Cancelling or killing request
7. Sandbox server upon client death
8. PID 1 process upon supervisor death
9. Signals (tracee, SIGPIPE, UBSan)
10. Running as superuser
11. Performance optimizations