# Intro

The idea of the work is to create minimalistic container to study techniques used in containerization engines and sandboxes. The parent process spawns child using `clone` function with namespaces isolation flags as parameter. Child process performs basic environment initialization: mounts of special filesystems like `proc` , `sys` , environmental variables, network interfaces configuration. Then, child process spawns shell with standard `fork-exec` sequence.

# Features

- Rootfs is saved from custom Docker image. It is based on the `debian:bullseye-slim` and contains `sysbench` utility for performance benchmarking, `ping` utility for network connectivity testing
- Host environment preparation is done via `configure.sh` bash script which performs extracting, setup, mount of rootfs and configuration of network interfaces to allow containers . This script is also responsible for cleaning artifacts left from the previous run
- PID isolation
    - As we can see on the screenshot, the number of processes within the container is much less than on average Linux distro. Moreover, there are only 3 numbered process entries in the `proc fs`

    

    

- Network namespace isolation & connection with host namespace
    - On the screenshot we can see possibility to ping host IP address from inside the container
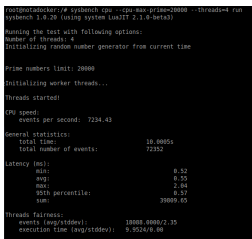
    

- Additionally, access to WAN might be granted to the container using NAT with the following command: `sudo iptables -t nat -A POSTROUTING -s 192.168.20.0/24 -o <outbound_interface> -j MASQUERADE` where `<outbound_interface>` is the interface on host from which traffic is going to the default gateway

# Run the project

```
./configure.sh
cmake -S . -B build
cd build && make && cd -
sudo build/ccontainer
```

# Performance comparison

- Note that before performing I/O tests, utility has to be created with `sysbench fileio --file-total-size=512M --file-num=1 prepare`
  - P.S. few files might be created instead of just one. It will affect the test and results might be slightly different

| Metric | Commands | CContainer | Docker | Host |
|---|---|---|---|---|
| CPU - 4 threads | `sysbench cpu --cpu-max-prime=20000 --threads=4 run` With multiple threads we can check whether (and how) multi threading affects performance in the container. | 10.0005s  | 10.0006s | 10.0003s |
| CPU - 1 thread | `sysbench cpu --cpu-max-prime=20000 --threads=1 run` Calculating max-prime number is CPU-bound task. Performance might be measured relative to the execution time. | 10.0003s | 10.0004s | 10.0001s |
| File IO Write | `sysbench fileio --file-total-size=512M --file-num=1 --file-test-mode=rndwr --file-rw-ratio=0 --time=30 run` This command performs benchmarks I/O with random writes. It outputs a number of various metrics such as number of writes performed. (non | written, MiB/s: 143.95 | written, MiB/s: 190.65 | written, MiB/s: 195.76 |

| Metric | Commands | CContainer | Docker | Host |
|---|---|---|---|---|
| | mixed with reads) However, we are interested in the throughput (MiB/S). | | | |
| File IO Read | `sysbench fileio --file-total-size=512M --file-num=1 --file-test-mode=rndrd --file-rw-ratio=0 --time=30 run` This command performs benchmarks I/O with random reads (non mixed with writes). This is artificial test which would approximately show us reads performance. Although, in real applications reads are often mixed with writes | read, MiB/s: 3326.45 | read, MiB/s: 3123.45 | read, MiB/s: 3254.85 |
| Memory | `sysbench memory --memory-block-size=4K --memory-total-size=2G run` Performs default memory write test with block size aligned to the page size. The command outputs the number of seconds needed to finish the task. It's important because it might show the ability of container to handle memory intensive tasks. | 1.8006s | 1.8962s | 1.8020s |

- For some reason, writes were slower in my container. I have went through possible causes variants, but unfortunately didn't find meaningful justification yet. The next step would be execution of container on different host (to ensure that the problem is not hardware / OS specific) and then trying lower level monitoring tools from the bcc collection.

# Sources

- https://linuxhint.com/use-sysbench-for-linux-performance-testing/
- https://github.com/89luca89/basic_c_container