



UNIVERSITATEA DIN BUCUREȘTI

FACULTATEA DE MATEMATICĂ ȘI INFORMATICĂ

Aplicație de gestionare a mediilor de virtualizare

LUCRARE DE LICENȚĂ

COORDONATOR ȘTIINȚIFIC

IROFTI PAUL

STUDENT

BUTNARU IOAN-SORIN

BUCUREȘTI, ROMÂNIA

IULIE 2020

Cuprins

I	Introducere	4
I.1	Prezentarea generală a temei	4
I.2	Scopul lucrării	5
I.3	Soluții existente	5
I.3.1	chroot	5
I.3.2	FreeBSD Jail	5
I.3.3	LXC	6
I.3.4	docker	6
I.4	Open Container Initiative (OCI)	6
I.5	Motivație	7
II	Tehnologii utilizate	8
II.1	Linux	8
II.1.1	namespace	8
II.1.2	cgroup	11
II.1.3	overlayFS	11
II.2	rust	11
III	Contribuția propriu-zisă a candidatului	12
III.1	Structurată în funcție de tipul lucrării, fiind compusă dintr-un număr arbitrar de capitole (vezi formatul lucrării în funcție de tip).	13
III.1.1	Descriu în detaliu fundamentarea teoretică și dezvoltarea aplicativă (dacă e cazul) a temei abordate	13
III.1.2	Conțin puncte de vedere personale, interpretări ale teoriilor și conceptelor abordate în lucrare	13
III.1.3	Trec în revistă abordări existente ale problemei cu evidențierea avantajelor și dezavantajelor	13

III.1.4	Descompun problema propusă de tema lucrării în subprobleme specifice și prezentarea modului de rezolvare, analize critice ale fenomenelor și proceselor studiate, comparații cu rezultate obținute anterior (unde e cazul), proiectarea aplicației, detalii de implementare, rezultate experimentale, exemple de test sau rezultate sub forma unor studii de caz, modul de utilizare a programului etc.	13
III.2	Arhitectura aplicației	13
III.2.1	pachete de clase pentru partea de user/admin	13
III.2.2	diagrama entitate relatie a bazei de date	13
III.2.3	alte componente incapsulate (maps, API-uri etc)	13
III.3	Funcționalitatea aplicației	13
III.3.1	un ghid al aplicației cu capturi de imagine și secvențe de cod reprezentative	13
III.4	Testarea aplicației	13
III.4.1	teste functionale	13
III.4.2	identificarea unor puncte slabe care vor fi rezolvate ulterior . . .	13
IV	Concluzii	14
IV.1	Soncluzii referitoare la modul de realizare a temei	14
IV.2	Posibile dezvoltări ulterioare ale temei	14
IV.3	Aprecierile personale privind relevanța rezultatelor obținute	14
IV.4	Eventuale aplicații practice sau interferențe cu alte domenii	14
IV.5	Descrierea aplicației, minimal (Am dezvoltat o aplicație care...)	14
IV.6	Directii viitoare de dezvoltare	14
	Bibliografie	15
	Anexe	17
A	Capitol anexă	18
A.1	Secțiune anexă	18

Abstract

Abstract-ul lucrării.

Capitolul I

Introducere

I.1 Prezentarea generală a temei

Containerizarea, sau virtualizarea la nivel de sistem de operare este o paradigmă a sistemului de operare în care *kernelul* permite existența a mai multor instanțe izolate ale utilizatorului. Astfel de instanțe, numite containere (Solaris, Docker), zone (Solaris), servere private virtuale (OpenVZ), partiții, medii virtuale (VEs), *kerneluri* virtuale (DragonFly BSD) sau închisori (FreeBSD jail sau chroot jail), pot părea ca niște computere reale din punctul de vedere al programelor care rulează în ele. Un program care rulează pe un sistem de operare obișnuit poate vedea toate resursele (dispozitive conectate, fișiere și foldere, *network share*-uri, putere a procesorului, capacități hardware cuantificabile) ale acelui computer. Cu toate acestea, programele care rulează în interiorul unui container pot vedea doar conținutul și dispozitivele destinate containerului. [1]

Pe sistemele de operare similare Unix, această caracteristică poate fi văzută ca o implementare avansată a mecanismului standard *chroot*, care schimbă folderul rădăcină aparent pentru procesul de rulare curent și copiii săi. În plus față de mecanismele de izolare, *kernelul* oferă și funcții de gestionare a resurselor pentru a limita impactul activităților unui container asupra altor containere. [1]

Virtualizarea la nivel de sistem de operare este folosită în mod obișnuit în mediile de găzduire virtuale, unde este utilă pentru alocarea în siguranță a resurselor hardware finite între un număr mare de utilizatori. Administratorii de sistem o pot utiliza, de asemenea, pentru consolidarea hardware-ului unui server, prin mutarea serviciilor de pe gazde separate pe un singur server în containerele. [1]

Alte scenarii tipice includ separarea mai multor programe în mai multe containere pentru o securitate mai bună, independență hardware și funcții adiționale de gestionare a resurselor. Implementările de virtualizare la nivel de sistem de operare, capabile să migreze direct, pot fi de asemenea utilizate pentru echilibrarea dinamică a resurselor utilizate de containere între nodurile dintr-un cluster. [1]

I.2 Scopul lucrării

Scopul acestei lucrări este de a dezvolta o aplicație de containerizare folosind funcționalitățile puse la dispoziție de către *kernelul* Linux. În dezvoltarea acestei aplicații, exemplele principale au fost aplicația Docker și specificațiile OCI.

I.3 Soluții existente

I.3.1 chroot

chroot este un apel de sistem care schimbă directorul rădăcină aparent pentru procesul de rulare curent și pentru copiii săi. Acesta a fost introdus în cursul dezvoltării versiunii 7 Unix în 1979 și a fost adăugat la BSD de Bill Joy la 18 martie 1982. [2]

Operația *chroot* nu este destinată să se apere împotriva manipulării intenționate de către utilizatorii privilegiați (root). În majoritatea sistemelor, contextele *chroot* nu sunt adăugate corect și programele *chroot*-ate cu privilegii suficiente pot efectua un al doilea *chroot* pentru a ieși din mediul respectiv. Pentru a atenua riscul acestei slăbiciuni în materie de securitate, programele trebuie să renunțe la privilegiile de administrator imediat după operația de *chroot*, sau ar trebui utilizate în schimb alte mecanisme - cum ar fi închisorile FreeBSD. Unele sisteme, cum ar fi FreeBSD, iau măsuri de precauție pentru a preveni un al doilea atac *chroot*. [3]

I.3.2 FreeBSD Jail

jail este o implementare a virtualizării la nivel de sistem FreeBSD, care permite administratorilor de sistem să partiționeze un sistem FreeBSD în mai multe mini-sisteme independente numite închisori, toate care au același kernel. [4] Acesta este implementat printr-un apel de sistem, *jail(2)*, precum și o utilitate de tip userland, *jail(8)*. [5]

I.3.3 LXC

LXC (Linux Containers) este o metodă de virtualizare la nivel de sistem de operare folosită pentru rularea mai multor sisteme Linux (containere) izolate pe o gazdă utilizând un singur *kernel* Linux. [6]

Kernelul Linux oferă funcționalitatea cgroups care permite limitarea și prioritizarea resurselor (procesor, memorie, I/O, rețea etc.) fără a fi necesară pornirea unor mașini virtuale și, de asemenea, funcționalitatea de izolare a spațiului de nume care permite izolarea completă a unei aplicații vederea a mediului de operare. [6]

LXC combină cgrupurile nucleului și suportul pentru spațiile de nume izolate pentru a oferi un mediu izolat pentru aplicații. Versiunile anterioare de Docker au folosit LXC ca driver de execuție a containerelor, deși LXC a devenit opțional în v0.9, iar suportul a fost oprit în Docker v1.10. [7][6]

I.3.4 docker

Docker este un set produse de tipul *platform as a service* (PaaS) care utilizează virtualizarea la nivel de sistem de operare pentru a livra software în pachete numite containere. Containerele sunt izolate unele de celelalte și conțin propriile programe, biblioteci și fișiere de configurare; ele pot comunica între ele prin canale bine definite. Toate containerele sunt administrate de un *kernelul* sistemului de operare și, prin urmare, folosesc mai puține resurse decât mașinile virtuale. [8]

I.4 Open Container Initiative (OCI)

Inițiativa Open Container (OCI) este proiect realizat de către fundația Linux, cu scopul de a crea standarde în jurul formatelor de containere și în jurul rulării acestora. OCI a fost lansată pe 22 iunie 2015 de Docker, CoreOS și alți lideri din industria containerelor.

OCI conține în prezent două specificații: Specificația pentru rulare (runtime-spec) și Specificația pentru imagini (image-spec). Specificația pentru rulare prezintă modul de a rula un „pachet de sisteme de fișiere” (*filesystem bundle*) care este despachetat pe disc. La un nivel înalt, o implementare OCI ar descărca o imagine OCI, apoi ar

despacheta imaginea în sistemul de fișiere al unui *runtime* OCI. Acest pachet urmează să fie rulat de către *runtime*.

Formatul de imagine OCI conține suficiente informații pentru a lansa o aplicație pe o platforma țintă (de exemplu, comandă, argumente, variabile de mediu etc.). Această specificație definește modul de creare a unei imagini OCI, care va fi realizată în general de către un sistem de compilare, și de generare a unui *image manifest*, a unei serializări de sistem de fișiere (*layer*) și a unei configurații a imaginii. La un nivel înalt, *image manifest*-ul conține metadate despre conținutul și dependențele imaginii. Configurația imaginii include informații, cum ar fi argumentele aplicației, medii, etc.

I.5 Motivație

Motivul pentru care am ales această temă este pentru a mă familiariza cu mecanismele care stau în spatele multor aplicații de virtualizare folosite în industrie și pentru a învăța un nou limbaj de programare, rust.

Capitolul II

Tehnologii utilizate

II.1 Linux

II.1.1 namespace

Un *namespace* înfășoară o resursă globală de sistem printr-o abstractizare care face ca procesele din *namespace* să creadă că au propria lor instanță izolată a resursei globale. Modificările aduse resursei globale sunt vizibile altor procese care sunt membre ale *namespace*-ului, dar sunt invizibile pentru alte procese.[9]

În prezent, Linux implementează opt tipuri diferite de *namespace*-uri.

Mount namespace

Mount namespace-ul (CLONE_NEWNS, Linux 2.4.19) izolează setul de *mount point*-uri a sistemelor de fișiere văzute de un grup de procese. Astfel, procesele din diferite *mount namespace*-uri pot avea vederi diferite ale ierarhiei sistemului de fișiere. Odată cu adăugarea *mount namespace*-urilor, apelurile de sistem *mount()* și *umount()* au încetat să funcționeze pe un set global de *mount point*-uri vizibile pentru toate procesele din sistem și în schimb au efectuat operațiuni care au afectat doar *mount namespace*-ul asociat procesului de apelare. [10]

O utilizare a *mount namespace*-urilor este crearea de medii care sunt similare cu închisorile chroot. Cu toate acestea, spre deosebire de apelul de sistem *chroot()*, *mount namespace*-urile sunt un instrument mai sigur și mai flexibil pentru această sarcină. Sunt posibile și alte utilizări mai sofisticate ale *mount namespace*-urilor. De exemplu, *mount namespace*-urile separate pot fi configurate într-o relație *master-slave*, astfel

încât evenimentele de montare să fie propagate automat de la un *namespace* la altul; acest lucru permite, de exemplu, unui hard disk care este montat într-un *namespace* să apară automat în alte *namespace*-uri. [10]

Mount namespace-ul a fost primul tip de *namespace* care a fost implementat pe Linux, apărut în 2002. Acest fapt explică numele destul de generic "NEWNS" (scurt pentru "new namespace"). [10]

UTS namespace

UTS namespace-urile (CLONE_NEWUTS, Linux 2.6.19) izolează două identificatoare de sistem, *nodename* și *domainname*, returnate de apelul de sistem *uname()*; acestea putând fi setate folosind apelurile de sistem *sethostname()* și *setdomainname()*. În contextul containerelor, funcția *UTS namespace*-urilor permite fiecărui container să aibă propriul *hostname* și propriul *NIS domain name*. Acest lucru poate fi util pentru scripturile de inițializare și configurare care își adaptează acțiunile pe baza acestor nume. Termenul *UTS* derivă din numele structurii transmise apelului de sistem *uname()*: struct *utsname*. Numele acestei structuri provine, la rândul său, de la *UNIX Time-sharing System*. [10]

IPC namespace

IPC namespace-urile (CLONE_NEWIPC, Linux 2.6.19) izolează anumite resurse de comunicare între procese (*IPC*; *interprocess communication*), și anume, obiecte *IPC System V* și (de la Linux 2.6.30) cozi *POSIX* de mesaje. Caracteristica comună a acestor mecanisme IPC este că obiectele IPC sunt identificate prin alte mecanisme decât *pathname*-urile sistemelor de fișiere. Fiecare *IPC namespace* are propriul set de identificatori *IPC System V* și propriul sistem de fișiere pentru cozile *POSIX* de mesaje. [10]

PID namespace

PID namespace-urile (CLONE_NEWPID, Linux 2.6.24) izolează spațiul numerelor de identificare ale proceselor. Cu alte cuvinte, procesele din *PID namespace*-uri diferite pot avea același PID. Unul dintre avantajele principale ale *PID namespace*-urilor este faptul că containerele pot fi migrate între gazde păstrând aceleași PID-uri pentru procesele din interiorul containerului. *PID namespace*-urile permit, de asemenea, fiecărui container să aibă propriul său *init* (PID 1), "strămoșul tuturor proceselor" care

gestionează diverse activități de inițializare a sistemului și reface procesele copil orfane atunci când se termină. [10]

Din punctul de vedere al unei anumite instanțe de *PID namespace*, un proces are două PID-uri: PID-ul din *namespace* și PID-ul în afara *namespace*-ului, din sistemul gazdă. Un proces poate vedea (de exemplu, să vizualizeze prin */proc/PID* sau să trimită semnale cu *kill()*) doar procesele conținute în *PID namespace*-ul său și *namespace*-urile aflate sub acel *PID namespace*. [10]

Network namespace

Network namespace-urile (*CLONE_NEWNET*, începute în Linux 2.4.19 2.6.24 și completate în mare parte în aproximativ Linux 2.6.29) asigură izolarea resurselor de sistem asociate cu rețelelor. Astfel, fiecare *network namespace* are propriile dispozitive de rețea, adrese IP, tabele de rutare IP, director */proc/net* și numere de port. [10]

Network namespace-urile din rețea fac ca containerele să fie utile dintr-o perspectivă de rețea: fiecare container poate avea propriul dispozitiv de rețea (virtual) și propriile aplicații care se leagă la al porturile *namespace*-ului; reguli de rutare adecvate în sistemul gazdă pot direcționa pachetele de rețea către dispozitivul de rețea asociat unui anumit container. Astfel, de exemplu, este posibil să existe mai multe servere web containerizate pe același sistem gazdă, fiecare server fiind legat la portul 80 în *namespace*. [10]

User namespace

User namespace-urile (*CLONE_NEWUSER*, începute în Linux 2.6.23 și completate în Linux 3.8) izolează spațiile numerelor de identificare ale utilizatorului și grupului. Cu alte cuvinte, ID-urile de utilizator și de grup ale unui proces pot fi diferite în interiorul și în afara unui *user namespace*. Cel mai interesant caz este faptul că un proces poate avea un ID de utilizator normal neprivilejat în afara unui *user namespace*, având în același timp un ID de utilizator 0 în interiorul *namespace*-ului. Acest lucru înseamnă că procesul are privilegii root complete pentru operațiunile din *user namespace*, dar nu este privilegiat pentru operațiuni în afara *namespace*-ului. [10]

Cgroup namespace

Cgroup namespace-urile (CLONE_NEWCGROUP, Linux 4.6) izoleaza vizualizarea cgrupurilor unui proces, care pot fi accesate prin `/proc/[pid]/cgroup` și `/proc/[pid]/mountinfo`. [11]

Fiecare *namespace* are propriul set de directoare cgroup rădăcină. Aceste directoare rădăcină sunt punctele de bază pentru locațiile relative afișate în din fișierul `/proc/[pid]/cgroup`. Când un proces creează un nou *cgroup namespace* folosind `clone(2)` sau `unshare(2)` cu *flag*-ul CLONE_NEWCGROUP, directoarele sale cgroup curente devin directoarele cgroup rădăcină ale noului *namespace*. [11]

Time namespace

Time namespace-urile izoleaza valorile a două ceasuri de sistem: CLOCK_MONOTONIC și CLOCK_BOOTTIME.[12]

Procesele dintr-un *time namespace* împărtășesc valorile pentru aceste ceasuri. Acest lucru afectează diferite API-uri care măsoară folosind aceste ceasuri, inclusiv: `ceas_gettime(2)`, `ceas_nanosleep(2)`, `nanosleep(2)`, `timer_settime(2)`, `timerfd_settime(2)` și `/proc/uptime`. [12]

În prezent, singura modalitate de a crea un *time namespace* apelând `unshare(2)` cu *flag*-ul CLONE_NEWTIME. Acest apel creează un nou *time namespace*, dar nu plasează procesul de apelare în noul *namespace*. În schimb, copiii ulterior creați ai procesului de apelare sunt plasați în noul *namespace*. Aceasta permite compensări de ceas pentru ca noul *namespace* să fie setat înainte ca primul proces să fie plasat în *namespace*. Legătura simbolică `/proc/[pid]/ns/time_for_children` arată *time namespace*-ul în care copiii procesului for fi creați. [12]

II.1.2 cgroup

II.1.3 overlayFS

II.2 rust

Capitolul III

Contribuția propriu-zisă a candidatului

- III.1 Structurată în funcție de tipul lucrării, fiind compusă dintr-un număr arbitrar de capitole (vezi formatul lucrării în funcție de tip).
 - III.1.1 Descriu în detaliu fundamentarea teoretică și dezvoltarea aplicativă (dacă e cazul) a temei abordate
 - III.1.2 Conțin puncte de vedere personale, interpretări ale teoriilor și conceptelor abordate în lucrare
 - III.1.3 Trec în revistă abordări existente ale problemei cu evidențierea avantajelor și dezavantajelor
 - III.1.4 Descompun problema propusă de tema lucrării în subprobleme specifice și prezentarea modului de rezolvare, analize critice ale fenomenelor și proceselor studiate, comparații cu rezultate obținute anterior (unde e cazul), proiectarea aplicației, detalii de implementare, rezultate experimentale, exemple de test sau rezultate sub forma unor studii de caz, modul de utilizare a programului etc.

III.2 Arhitectura aplicației

III.2.1 Arhitectura de înalt nivel a aplicației / subaplicației

Capitolul IV

Concluzii

- IV.1 Soncluzii referitoare la modul de realizare a temei
- IV.2 Posibile dezvoltări ulterioare ale temei
- IV.3 Aprecierile personale privind relevanța rezultatelor obținute
- IV.4 Eventuale aplicații practice sau interferențe cu alte domenii
- IV.5 Descrierea aplicatiei, minimal (Am dezvoltat o aplicatie care...)
- IV.6 Directii viitoare de dezvoltare

Bibliografie

- [1] Wikipedia contributors. Os-level virtualization — Wikipedia, the free encyclopedia. https://en.wikipedia.org/w/index.php?title=OS-level_virtualization&oldid=958577453, 2020. [Online; accessed 7-June-2020].
- [2] Poul-Henning Kamp and Robert NM Watson. Jails: Confining the omnipotent root. In *Proceedings of the 2nd International SANE Conference*, volume 43, page 116, 2000.
- [3] chroot - freebsd manual pages. <http://man.freebsd.org/chroot/2>. [Online; accessed 7-June-2020].
- [4] David Chisnall. Dragonfly bsd: Unix for clusters? <https://www.informit.com/articles/prINTERfriendly/766375>, 2020. [Online; accessed 2019-03-06].
- [5] Wikipedia contributors. FreeBSD jail — Wikipedia, the free encyclopedia. https://en.wikipedia.org/w/index.php?title=FreeBSD_jail&oldid=945967483, 2020. [Online; accessed 7-June-2020].
- [6] Wikipedia contributors. Lxc — Wikipedia, the free encyclopedia. <https://en.wikipedia.org/w/index.php?title=LXC&oldid=958749686>, 2020. [Online; accessed 11-June-2020].
- [7] Docker 0.9: introducing execution drivers and libcontainer - Docker Blog. <https://www.docker.com/blog/docker-0-9-introducing-execution-drivers-and-libcontainer/>. [Online; accessed 2018-05-09].
- [8] Wikipedia contributors. Docker (software) — Wikipedia, the free encyclopedia. [https://en.wikipedia.org/w/index.php?title=Docker_\(software\)&oldid=960398220](https://en.wikipedia.org/w/index.php?title=Docker_(software)&oldid=960398220), 2020. [Online; accessed 7-June-2020].

- [9] namespaces(7) - linux programmer's manual. <https://man7.org/linux/man-pages/man7/namespaces.7.html>. [Online; accessed 11-Jun-2020].
- [10] Namespaces in operation, part 1: namespaces overview. <https://lwn.net/Articles/531114/>. [Online; accessed 11-Jun-2020].
- [11] cgroup_namespaces(7) - linux programmer's manual. https://man7.org/linux/man-pages/man7/cgroup_namespaces.7.html. [Online; accessed 11-Jun-2020].
- [12] time_namespaces(7) - linux programmer's manual. https://man7.org/linux/man-pages/man7/time_namespaces.7.html. [Online; accessed 11-Jun-2020].

Anexe

Anexa A

Capitol anexă

A.1 Secțiune anexă

Conținut anexă