

異種 OS 混在環境におけるコンテナ型仮想化のライブマイグレーション

b1014223 高川雄平
指導教員：松原克弥

Container Live Migration among Heterogeneous OS Platforms

Yuhei Takagawa

概要：オペレーティングシステム (OS) が提供する計算資源や名前空間をアプリケーション毎に多重化して隔離することで、ひとつの OS 上に独立した複数の実行環境を実現するコンテナ型仮想化が注目されている。一方、ハードウェア資源を仮想化する従来方式と比較して、コンテナ型仮想化の実現方式は OS への依存度が高く、異なる OS 上で実現されたコンテナ環境間での互換性がない。本研究では、クラウドコンピューティング環境において、負荷分散や可用性実現の要として広く利用されているライブマイグレーション技術に着目して、動作中のコンテナ実行環境を異なる種類の OS が動作する別の計算機へ動的に移動して実行の継続を可能にする方式の検討と実装を行う。OS 毎のシステムコールや ABI の変換、プロセス実行状態の取得と復元方式の統一化、資源隔離機能の対応づけにより、Linux と FreeBSD を対象とした異種 OS 間のライブマイグレーションでの実現を目的とする。

キーワード：コンテナ仮想化, ライブマイグレーション, FreeBSD, Linux

Abstract: The container-based virtualization, that multiplexes and isolates computing resource and name space which operating system (OS) provides for each application, has been recently attracted. As compared with the conventional hardware virtualization, containers run on different OSs may not be compatible because their container implementations must depend on each underlying OS. This research focus on the live migration since it is one of the most important technology for realizing load balancing and availability in cloud computing, that is a major application of the virtualization. This study clarifies how system calls and ABI can be converted, how running containers can be captured with an uniform representation and they restore in both OSs, and how an uniform resource isolation can be realized, especially for Linux and FreeBSD as the 1st target.

Keywords: the container-based virtualization, live migration, FreeBSD, Linux

1 背景と目的

クラウドコンピューティングを支える技術として、仮想化技術がある。特に、アプリケーション毎の実行環境の軽量なコンパートメント化・ポータビリティによる簡易的な環境構築を目的としたコンテナ型仮想化が注目されている [1]。コンテナ型仮想化は、OS が提供する資源を分離・制限し、単一動作の OS 上に複数の独立した実行環境を構築できる軽量な仮想化を実現する。

また、クラウドコンピューティング実動環境では、負荷分散と可用性を実現するためにライブマイグレーションが活用されている。ライブマイグレーションは、実行中のサービスを動的に別のマシンに移行する技術である。コンテナ型仮想化におけるライブマイグレーションも実現されており、Linux では CRIU[2, 3], FreeBSD では FreeBSD VPS(Virtual Private System)[4]

が実装として存在する。OS 依存度が大きいコンテナ型仮想化では、OS の計算資源や API, コンテナ実現方式が異なるため、異なる OS 環境間でコンテナをマイグレーションすることができない。同一 OS 間のライブマイグレーションでは、システム全体の計算資源を活用することができず、メンテナンス時の可用性が損なわれてしまう。

本研究では、稼働 OS に依存したシステム運用で起きる前述の問題に対処するために、OS 混在環境におけるコンテナ型仮想化実行環境のライブマイグレーションの実現を目的とする。

2 関連技術

2.1 コンテナ型仮想化

コンテナ型仮想化は、アプリケーションが使用する計算資源やアクセス可能な名前空間を分離することでアプリケーションごとに異なる実

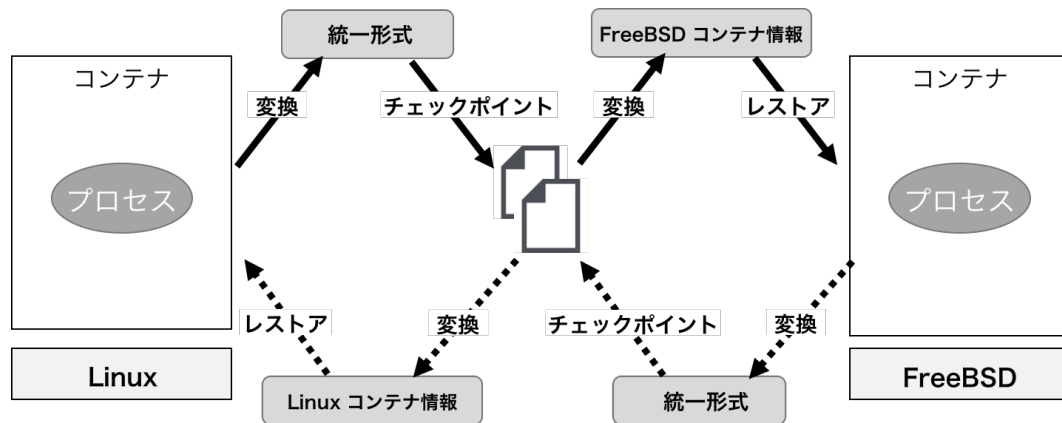


図 1 提案システムの概要

行環境を構築できる技術である。コンテナとは分離されたアプリケーション実行環境のことを指す。コンテナの実現には、LinuxではcgroupsやNamespace, FreeBSDではJailというOS機能が利用されている。

cgroupsは、Linuxのプロセスが使用する資源を制限できる機能である。例えば、CPU使用率やメモリ利用量などを制限可能である。Namespaceは、LinuxのプロセスIDやディレクトリ、ネットワークなどを制御するための名前空間を分離できる機能である。

JailはLinuxにおけるNamespaceにあたる機能で、ディレクトリやネットワークなどの識別子を分離する機能である。分離した空間は他の空間には直接関係しないため、同じ識別子を分離した空間ごとに利用することができる。

2.2 ライブマイグレーション

ライブマイグレーションは、仮想化環境上で動作しているプロセスを停止させずに仮想化環境ごと別のマシンに移動する技術で、負荷分散や可用性の実現を目的に用いられる。ライブマイグレーションには3つの手順がある。最初に、CPUやメモリなどのハードウェア資源(以下、H/W資源)の状態をファイルに保存する。次に保存したファイルを転送する。最後に、転送されたファイルからH/W資源の状態を復元する。H/W資源の状態を保存するための処理を「チェックポイント」と呼び、H/W資源の状態を復元することを「レストア」と呼ぶ。

Linuxでライブマイグレーションを実現可能なCRIU, FreeBSDでライブマイグレーションを実現したFreeBSD VPSについて以下に述べる。

2.3 CRIU

CRIU(Checkpoint/Restore in Userspace)は、Linux上でコンテナのマイグレーションを可能にするOSS実装である。CRIUの機能は、主にプロセス実行・隔離状態のチェックポイントを作成することとプロセス実行・隔離状態のレストアを行うことである。ライブマイグレーションは状態が移動するだけなので、マイグレーション先のマシンには同じコンテナが存在する必要がある。CPUの復元では、取得したレジスタの値をptraceシステムコールの機能を使って設定することでCPUの状態を復元している。なお、チェックポイント/レストアに必要なユーザ空間では行えない処理をLinuxカーネルに組み込んでおり、他の数多くの機能はユーザ空間で実現しているという特徴がある。

2.4 FreeBSD VPS

FreeBSD VPS(Virtual Private System)は、FreeBSD jailのマイグレーションを可能にしたシステムであり、VPSインスタンスというコンテナを作成する。CRIUはプロセスをマイグレーションしているが、FreeBSD VPSはVPSインスタンスごとマイグレーションしている。CPUの復元にはPCB(Process Control Block)を活用しているため、CRIUと全く異なる方法でプロセスの復元を行っている。

3 異種OS間ライブマイグレーション実現方式

本章では、異種OS間でのコンテナ型仮想化におけるライブマイグレーションの実現方式を提案する。本研究では、LinuxとFreeBSDを対

象とする。図1は、LinuxとFreeBSD間におけるコンテナ型仮想化のライブマイグレーションの概要図である。本提案では、プロセス実行・隔離状態などの情報を統一形式に変換してチェックポイント時にファイルに保存し、レストア時にファイルから取得した情報をOSに適した形式に変換してプロセス実行・隔離状態を復元する。本研究では、統一形式をLinuxでCRIUを用いて復元できる形式とする。

コンテナ型仮想化実行環境のマイグレーションのためには、プロセス実行状態のマイグレーション(プロセス・マイグレーション)とプロセス隔離状態のマイグレーション(コンテナ・マイグレーション)が必要になる。それぞれにおける技術的課題と解決方法を第3.1章、第3.2章に示す。

3.1 プロセス・マイグレーションの実現

プロセス・マイグレーションを実現するには、以下の2点の課題を解決する必要がある。

3.1.1 システムコールの差異

LinuxとFreeBSDの共通なシステムコールの大半は同じ動作である。しかし、システムコール番号やシステムコールの引数・パラメータの値が異なる。そのため、LinuxバイナリはFreeBSD上で動作せず、FreeBSDバイナリはLinux上で動作しない。

また、システムコール引数の渡し方がLinuxとFreeBSDでは異なる。FreeBSD/x86はシステムコール引数をスタック経由で渡すのに対して、Linux/x86はシステムコール引数をレジスタ経由で渡す(表1参照)。Linux/x86では、システムコールの引数の上限が5つに決まっており、それ以上は利用できない。関数の引数はスタック経由であるため、6つ以上の引数を持つことができる。そのため、同じプログラムが動作していても、FreeBSDとLinuxではメモリやレジスタの状態が異なる。

システムコールの差異に関しては、システムコールの変換を行うことで解決する。システムコール番号や引数・パラメータ値に関しては、それぞれを変換することでOSが異なってもシステムコールを動作させることができる。システムコール引数の渡し方に関しては、FreeBSD上でLinuxバイナリ実行時には引数をレジスタに入れ、Linux上でFreeBSDバイナリ実行時には

引数をスタックに積むというような実装を行うと、異なるOSでもレジスタやメモリの状態を同じにすることができる。

FreeBSDには、Linuxバイナリ互換機能[5]というカーネル機能があり、システムコール番号や引数・パラメータの変換・引数の渡し方の互換を行っている。この機能によって、Linuxバイナリであれば、LinuxとFreeBSDで実行ファイルを一切変更せずに動作させることができる。本提案では、システムコールの差異の吸収をLinuxバイナリ互換機能を用いることで実現する。

3.1.2 メモリレイアウトの差異

メモリレイアウトは仮想メモリ空間における領域の配置である。一般的にASLR(Address Space Layout Randomization)によって、領域の配置アドレスはランダムに割り当てられる。また、ASLRを無効化して領域の配置アドレスを固定した場合でも、FreeBSDとLinuxではスタック領域に0x1000番地の誤差がある。この誤差を純粋に修正する場合には、スタック内にある全てのベースアドレス、リターンアドレスなどを見つけ出し0x1000番地移動する必要がある。

メモリレイアウトの差異に関しては、メモリレイアウト変更機能を実装することで解決する。Linuxでは、prctlシステムコール機能を利用することでメモリレイアウトをユーザプログラムから変更することができる。そのため、チェックポイントのメモリマップを変更せずに、0x1000番地の誤差をなくすことができる。また、ASLRのランダムなメモリマップも再現可能である。しかし、FreeBSDにはメモリレイアウトをユーザプログラムから変更できる機能はない。LinuxとFreeBSDとの相互マイグレーションやASLRに対応してマイグレーションを実現するには、FreeBSD上でメモリレイアウトを変更できるシステムコールを作成する必要がある。

3.2 コンテナ・マイグレーションの実現

第2.1章で述べたように、FreeBSDとLinuxではコンテナを実現する機能が異なる。隔離・制限する対象や名称、値が異なるため、チェックポイントの状態をそのまま復元することはできない。

隔離・制限機能の差異に関しては、FreeBSDのJailとLinuxのcgroup、Namespaceの機能の対

表 1 x86 における Linux と FreeBSD のシステムコール引数格納場所

引数	Linux	FreeBSD
第 1 引数	EBX	スタック
第 2 引数	ECX	
第 3 引数	EDX	
第 4 引数	ESI	
第 5 引数	EDI	
第 6 引数	×	

表 2 プロセス制限機能対応

制限機能	Linux	FreeBSD
CPU リソース	cgroups	RCTL+cpuset(1)
メモリリソース		RCTL
ファイル入出力		
デバイスアクセス		devfs(8)
一時停止/再開		△ SIGSTOP

表 3 プロセス隔離機能対応

隔離機能	Linux	FreeBSD
プロセス間通信	Namespace	jail
マウントポイント		
ユーザ ID		
ホスト名		
ネットワーク		
プロセス ID		△ jail

応付けや数値の変換を行い、それぞれの機能で隔離・制限状態を再現することで解決する。資源の制限に関しては、FreeBSD の RCTL と Linux の cgroups の制限状態の対応付けと相互変換 (図 2 参照) を行うことで再現する。例えば、CPU 使用率を 50% に制限するのであれば、Linux 上では cgroups に以下を設定し、
`cpu.cfs_quota=50`
`cpu.cfs_period_us=100`
FreeBSD 上では RCTL に以下を設定することで
`jail:<jail id>:pcpu:deny=50`
同じ制限にすることができる。隔離状態の対応付けに関しても図 3 の対応付けを行うことで、隔離状態を再現できる。

4 実装状況

FreeBSD で FreeBSD VPS を用いない単純なプロセス・マイグレーションを実現した。ptrace() を用いてレジスタの状態を取得し、procfs の mem

からプロセスのユーザ空間のメモリを取得する。復元時には、プロセスをフォークしプログラムを実行開始前で停止させ、レジスタやメモリの状態を書き込むことで、計算処理の続きからプロセスを再開することができる。ただし、プロセスを動作させるために最低限必要な情報しかないため、ネットワークやファイル書き込みの状態を復元することはできない。また、FreeBSD から Linux へ単縦なプロセス・マイグレーションを実装している。Linux バイナリを対象として、Linux バイナリ互換機能を活用することでシステムコールの差異を変換し、Linux の prctl() を利用することでメモリレイアウトの差異を吸収することで実現する。

5 おわりに

本研究は、異種 OS 間でのコンテナ型仮想化のライブマイグレーションの実現を目的としている。対象 OS を Linux と FreeBSD とし、プロセス・マイグレーションとコンテナ・マイグレーションの検討を行った。プロセス・マイグレーションでは、システムコールの差異を Linux バイナリ互換機能を用い、メモリレイアウトの差異を prctl システムコール機能を用いる。提案するシステムを用いると、稼働 OS に依存しないシステム運用が可能になり、計算資源の活用や可用性の向上を見込める。

今後の課題は、コンテナ・マイグレーションの実現、ネットワークやファイルを利用するプロセスへの対応などが挙げられる。

参考文献

- [1] 451 Research, "451 Research: Application containers will be a \$2.7bn market by 2020", 2017.
- [2] CRIU Project: CRIU Main page, <https://criu.org/>, 2010 (accessed November 1 2017).
- [3] A. Mirkin, A. Kuznetsov and K. Kolyshkin: Containers checkpointing and live migration, In Proceedings of the 2008 Ottawa Linux Symposium, Vol. 2, 2008, pp. 8590.
- [4] Klaus P. Ohrhallinger: Virtual Private System for FreeBSD, <http://www.7he.at/freebsd/vps/>, 2010 (accessed November 6, 2017). EuroBSD-Con 2010.
- [5] Jim Mock: Linux Binary Compatibility, <https://www.freebsd.org/doc/handbook/linuxemu.html>, 1999 (accessed November 2, 2017).