





目次

第 1 章	はじめに	1
第 2 章	KVM	3
2.1	QEMU	3
2.2	KVM	4
第 3 章	Containers 超入門	7
3.1	Containers の世界と LXC、そして Docker	7
3.2	LXC を使い軽量仮想環境を手に入れよう	10
第 4 章	あとがき	15
4.1	こじろー	15
4.2	まっきー	15
4.3	だーまり	15



第1章

はじめに



第2章

KVM

KVM について知りたい人がいたとして、その人が知りたいのは恐らく KVM ではなく QEMU だろう。この章では、実は KVM は、特に何もしていない、というのを解説しようと思う。メモリの割り当て？ネットワーク接続？USB？ディスク？それは KVM ではなく QEMU で語るものだ。皆が何となく「ハイパーバイザ」と呼ぶ物、それは正確には「仮想マシン」と呼び、それは QEMU によって実現される。KVM は QEMU の仮想 CPU 高速化モジュールとしてとらえてしまっても、特に誤解はない。そして、これを正確には「ハイパーバイザ」と呼ぶ。

KVM と QEMU の違いをはっきりさせないと、有らぬ誤解を生むばかりか、適切な情報を検索することもままならない。この2つは似て非なるものではなく、はっきり言って、特に依存はない。KVM が無くても QEMU は動く。実は KVM なんて影も形もない頃から、QEMU は比較的安定して動作していた。あえて言ってしまうと、むしろ KVM の方が、存在意義という観点から、QEMU に依存している。

勉強し始めた時、僕はまず KVM から勉強を始めた。そこから vhost や virtio をいったものを調べ始めた。が、どうしてもよく分からない。最近ようやく理解のとっかかりをつかみ始めたが、ここに至ってようやく、最初から QEMU を勉強しておけば迷いにくいことに気付いた。というわけで、この本では QEMU の解説からしていくことにする。

2.1 QEMU

フォン・ノイマンという人が提唱したノイマン型コンピュータの仕組みは、何の改善もないまま、現在のコンピュータに使われている。基本的には、こうだ。コンピュータは以下の要素で成り立つ

- 演算装置 (CPU)
- 記憶装置 (RAM, Disk)
- 出力装置 (Tape, Display)
- 入力装置 (マウス・キーボード)

記憶装置に書き込まれたプログラム通りに、記憶装置に一時的な計算結果を書き込みつつ、演算装置が計算を行う。可変パラメータは入力装置から得られるし、計算結果は出力装置に送られる。言ってしまうと、コンピュータはこれだけのことしかない。CPU にはレジスタもある、という反論もあるが、それは要するに記憶装置だ。ネットは？それはよく訓練された入出力装置に過ぎない。まとめてしまえばたったこれだけでコンピュータは動作する。

そして、驚くべきことに、QEMU を知る上での前提知識は以上だ。QEMU とは、つまりこの各装置をソフトウェアでエミュレートしているだけなのだ。

2.1.1 CPUのエミュレート

KVM の話をする時に必要になるのが、CPU の仮想化だ。CPU が解釈できる様々な演算命令をすべてエミュレートすれば、CPU の仮想化ができる。QEMU は様々な CPU の命令をエミュレートできる。ARM も x86 も認識できる。

中でやっていることは単純だ。実行すべき命令をホストの CPU 命令に書き換えて実行する。x86 の CPU で動作するホストで ARM の命令をエミュレートするには、ARM の命令を x86 の命令に変換すればよい。CPU でできることなどたかがしれているから、ほぼ一対一で対応がある。対応がない場合は、複数の命令を組み合わせで代替する。

で、問題は x86 の命令をエミュレートして x86 の命令に書き換える時だ。これはエミュレートする必要がない。直接ホストの x86 の命令を実行すればよい。双方同じ CPU モデルなのだから、同じ命令がそのまま実行できるはずだ。そっちの方が変換の手間がなくなるから実行が高速になるが、それ以前に、そもそも無駄だ。

のだが、ここで問題が起きる。実は、そのまま実行できない命令があるのだ。これは CPU が悪いのでも、QEMU のエミュレートが悪いわけでもない。実行できないのは、OS が実行を制限してしまうからだ。

2.1.2 リングプロテクション

実はすべての CPU 命令をユーザーが実行できるわけではない。OS しか実行できないような強力なコマンドは、OS のユーザーは実行することができない。これを OS の「リングプロテクション」と言い、これら特別な命令群を「Ring 0」とか「センシティブな命令」と読んだりするが、要するに OS のような特別な奴にしか実行が許されていない命令群があるのだ。Linux で言うならば、カーネル空間でしか実行できない命令群である。

ここで常に念頭において置かねばならないが、当たり前で忘れがちなことがある。それは、エミュレーターである QEMU は「ユーザー空間のプログラムである」ということである。QEMU がユーザー空間で動作するソフトウェアであることが、脈々と発展し続ける KVM 開発の原動力となる。すべての周辺技術はこの当たり前の事実との戦いの歴史である。

つまり、確かに QEMU は x86 の命令をすべて解釈し実行できる能力はあるが、センシティブな命令は QEMU が動作している OS により実行が禁止される。つまり、実行できない。実行しようとした場合はどうなるかわからない。普通なら、強制的に Kill されるだろう。ホストから見れば、完全に「悪意のあるプログラム」にしか見えないからだ。

センシティブな命令を実行するには、システムコールを使うか別のセンシティブでない命令群で置き換えるかしなければならない。さて、今 QEMU は命令をホストの CPU で直接実行している。センシティブな命令が実行されようとしていることを知るにはどうすれば？この瞬間を QEMU がキャッチできなければ、QEMU はプロセスごと Kill されることになる。

長かったが、ここで KVM が現れる準備が整ったことになる。

2.2 KVM

さて、前章をもってすれば KVM の解説は単純を極める。KVM とは、センシティブな命令が実行されようとしていることを、ユーザー空間のプログラムに教える機能を持つ、カーネルモジュールである。知りたいユーザー空間のプログラムは `/dev/kvm` をポーリングしていると、KVM が教えてくれる。KVM はそれしかしない。今の KVM はいかに「本当にそれだけのことをする」ように機能を削り、磨かれているところである。

2.2.1 OpenStack

- KVM
- OpenStack

1. KVM
2. Docker
3. LXC

```
./stack.sh
```



第3章

Containers 超入門

3.1 Containers の世界と LXC、そして Docker

3.1.1 昔からあるコンテナ技術

コンテナ技術を取り囲む現在の状況と、それを踏まえた上での LXC と Docker の根本的な違いについて説明したいと思う。Linux Containers(LXC) は、どうやって我々がアプリケーションを動かすスケールさせるかという問題を変化させる可能性を持っている。コンテナ技術は新しいものではない。そして、LXC に関して言うと、追加パッチを Linux Kernel に適用させることなく、vanilla Linux Kernel 上で稼働させることができる。なお、LXC の Version1 は、長期サポートバージョンであり、5 年間サポートされることになる。話が逸れるが、vanilla Linux Kernel とは、Linux 作者の Linus Torvalds 氏がリリースするプレーンな Kernel のことである。それをベースに様々なベンダーが追加で拡張していくのである。また、vanilla という言葉には「普通の、ありきたりな、おもしろみのない」という意味がある。

話を戻そう。

コンテナ技術は、最近登場した新技術ではない。昔から存在し色んな所で採用されている。FreeBSD には Jail があり、Solaris には Zone がある。それに加えて、OpenVZ や Linux VServer のような Containers も存在する。その歴史は、chroot に始まり、FreeBSD Jail を経て、Linux Containers(LXC) に至る。chroot では、大雑把に言って、ディレクトリツリーの分離を行っていた。プロセスリスト自体は共有するようなモデルである。chroot のユースケースとしては、開発者向けのテスト/ビルド用環境である。FreeBSD Jail では、chroot の機能に加えて、プロセスリストとネットワークスタックも分離 (というか隔離) された。ユースケースとしては、root 権限の一般ユーザへの委譲、またそれに頼る形でのホスティングサービスである。LXC では、リソース管理テーブルを隔離し、cgroups によるシステムリソース (CPU、メモリ、ディスク etc) の制御を行えるようになった。これにより、LXC は、軽量な仮想環境と見なすことができるようになった。

3.1.2 なぜ皆コンテナに騒いでいるのか

コンテナは、ホストシステムからアプリケーションのワークロードを隔離、あるいはカプセル化する。コンテナを、ホスト OS 内にあるアプリケーションが実行されている OS と見なすことができ、かつ、それは Virtual Machine のように振る舞うのである。このエミュレーションは、Linux Kernel それ自体と、様々なディストリビューションとコンテナを使ってアプリケーションを動かすユーザのためにコンテナ用 OS のテンプレートを提供する LXC Project によって、実現されている。このように、Containers 技術が仮想マシンのように振る舞うことが可能になったことが、一気に注目を浴びる原因となったのである。

3.1.3 コンテナの価値は？

コンテナは、アプリケーションをホスト OS から分離、抽象化することで、LXC をまたぐシステム間でのポータビリティをもたらす。また、コンテナは、ハードウェアをエミュレートすることなく、cgroups と namespaces を駆使して Linux Kernel 内でベアメタルに近いスピードの軽量な OS 環境を実現している。シンプルで高速、かつハードウェアの仮想化よりもよりポータビリティがありスケールしやすいという構造により、コンテナは、根本的なユーザのワークロードやアプリケーションの仮想化の方法を変えるものなのである。なお、ここでいうポータビリティとは Docker によりもたらされる、どこでも同一のアプリケーションを稼働することができるという意味ではない。

3.1.4 LXC

LXC プロジェクトは、コンテナ用の OS テンプレートとライフサイクル管理のための幅広いツールセットを提供しています。現在、Canonical のサポートのもと、Stephane Graber と Serge Hallyn により開発は主導されています。

LXC は、活発に開発されているがその割にはドキュメントが少ない。特に Ubuntu 以外のディストリビューションで利用する際のドキュメントが欠如しており、多くの機能はまず Ubuntu 上で実装される。他のディストリビューションを利用しているユーザーからしたら、とてもフラストレーションのたまることである。また、ネット上には数多くの誤解を招くような情報があふれ、混乱を招いている状況も少なからずある。広くマーケットで存在感を示している Docker と混同されたり、そもそもの情報の多さが混乱の元となっていたりする。

LXC は、Docker のようなフロントエンドのアプリケーションのためのローレイヤー層なのか、はたまた、Docker が LXC 上に構築されたユーザーフレンドリーなフロントエンドなのか？こういった不確かな情報が広く出回っている。コンテナ技術のメリットを享受するために必ずしも Docker を使う必要はない。Docker はコンテナ利用の一つの選択でしかないのである。

3.1.5 Docker と LXC では何が違うのか

Docker 視点から見た LXC との違いについて説明する。そもそも LXC に対して Docker が提供している機能とは何なのか。まず第一に言われることは、どのようなホスト OS であってもポータブルなデプロイが可能である点である。Docker はアプリケーションをビルドするためのフォーマットが定義されている。この定義をまとめて記述するために Dockerfile ファイルを利用する。この Dockerfile ファイルはビルドでよく使われる makefile ファイル同様に Docker Containers の構成情報をまとめて記述するテキストファイルである。このファイルに記述する定義情報が全ての依存関係をカプセル化しているため、それはどこで実行してもアプリケー

ション実行環境が同一になるのです。LXC のプロセスのサンドボックスもポータビリティを持っているが、もし LXC のコンフィグをカスタマイズしているとしたら、ネットワークやストレージ、ディストリビューションの違いにより、それは別環境で稼働しない可能性が高くなります。Docker はその全てを抽象化するためどんな環境でも稼働させることができるのである。

Docker について言及するエンジニアは、総じてアプリケーション寄りのエンジニアである。Docker は、軽量の仮想マシンとしての利用というよりもアプリケーションのデプロイに最適化されている。これは Docker 自体の API やデザインの設計思想に反映されている。それとは対照的に、LXC は軽量の仮想マシンとしての利用に注力している。Docker には、git に似たバージョン管理機能が含まれている。バージョン間の diff の取得や Commit、ロールバックが可能となっている。それによって、Containers の変更を誰がどのように行ったのかについての全てのログを追うことができる。

他にも Docker の利点は多く存在するが、主にこのような点により、Docker は、コンテナそのものに対する見方を変えるきっかけを作った。今まで軽量の仮想マシンとして見られていたコンテナ技術をアプリケーションとしてのコンテナとしてエンジニアに再認識させることに成功したのである。

3.2 LXC を使い軽量仮想環境を手に入れよう

LXC の基本的なコマンドを使ったコンテナ操作を、Ubuntu14.04 をベースにした環境を使って説明していきたい。

3.2.1 LXC のインストール

Ubuntu の最新版である Ubuntu 14.04 LTS では、LXC 1.0.7 が `lxc` というパッケージ名で提供されている。また、Debian 8 (Jessie) では、LXC 1.0.6 のパッケージが提供されている。インストールは以下のコマンドを叩くだけである。

```
$ sudo apt-get install lxc
```

3.2.2 LXC で仮想環境を立ち上げる

LXC による仮想環境を立ち上げるためには、まずテンプレートと呼ばれる設定ファイルを用いる。デフォルトでメジャーディストリビューションのテンプレートはすでに同梱されているためこちらを利用する。テンプレートは `/usr/share/lxc/templates/` に配置されている。

```
$ ls /usr/share/lxc/templates/  
lxc-alpine      lxc-archlinux  lxc-centos     lxc-debian     lxc-fedora     lxc-openmandriva  
lxc-oracle      lxc-sshd       lxc-ubuntu-cloud  
lxc-altlinux    lxc-busybox    lxc-cirros     lxc-download   lxc-gentoo     lxc-opensuse  
lxc-plamo       lxc-ubuntu
```

Ubuntu のテンプレートを用いて `test-container-101` という名前の Ubuntu のコンテナを立ち上げる。

```
$ sudo lxc-create -t ubuntu -n test-container-101
```

これでコンテナの `root` ディレクトリに相当するディレクトリに必要なものがインストールされる。コンテナの場所は以下のディレクトリである。

```
/var/lib/lxcコンテナ名/<>/
```

そこで、`test-container-101` の `rootfs` の中を覗いてみると以下の通りとなる。

```
$ sudo ls -F /var/lib/lxc/test-container-101/rootfs/  
bin/  boot/  dev/  etc/  home/  lib/  lib64/  media/  mnt/  
opt/  proc/  root/  run/  sbin/  srv/  sys/  tmp/  usr/  var/
```

インストールしたコンテナの起動には以下のコマンドを実行する。

```
$ sudo lxc-start -n test-container-101 -d
```

-d オプションでデーモンとしてコンテナを起動する。この状態で lxc-console コマンドを用いてコンソール接続することでコンテナの内部にログインすることができる。

```
$ sudo lxc-console -n test-container-101
```

コンテナから抜ける際には、Ctrl+A を入力してその後 Q を押す。また、デフォルトのユーザは ubuntu で、パスワードも ubuntu である。インストール時に以下のメッセージが表示されているはずである。

```
# The default user is 'ubuntu' with password 'ubuntu'!  
# Use the 'sudo' command to run tasks as root in the container.
```

コンテナの終了は、lxc-shutdown コマンドを実行すればよい。

```
$ sudo lxc-shutdown -n test-container-101
```

3.2.3 コンテナの情報を見る

コンテナに関する情報を見てみよう。lxc-ls というコマンドでホスト上にあるコンテナの情報を確認することができる。-fancy オプションを付けることで、コンテナ名、状態、IPv4 のアドレス、IPv6 のアドレス、自動起動の有無を確認できる。

```
$ sudo lxc-ls --fancy  
NAME                STATE    IPV4    IPV6    AUTOSTART  
-----  
test-container-101  STOPPED -       -       NO
```

コンテナ単体の詳細情報については lxc-info というコマンドが提供されている。コンテナが STOPPED した状態のときにはこう表示される。

```
$ sudo lxc-info -n test-container-101  
Name:                test-container-101  
State:                STOPPED
```

コンテナを起動すると詳細情報が表示される。

```
$ sudo lxc-info -n test-container-101  
Name:                test-container-101  
State:                RUNNING  
PID:                 20434  
CPU use:              0.77 seconds  
BlkIO use:            7.16 MiB  
Memory use:           13.53 MiB  
KMem use:             0 bytes  
Link:                 vethABI04E  
TX bytes:             940 bytes  
RX bytes:             592 bytes  
Total bytes:          1.50 KiB
```

このように特定の情報のみを取得することも可能である。

```
$ sudo lxc-info -n test-container-101 -c lxc.utsname -c lxc.rootfs  
lxc.utsname = test-container-101  
lxc.rootfs = /var/lib/lxc/test-container-101/rootfs
```

3.2.4 LXD とは何か

LXD について本家ページを元に少し説明する。LXD とは Linux Container Daemon の略である。Canonical 主導で開発が進められているコンテナ技術であり、コンテナに今どきのハイパーバイザーの機能を追加するサーバプログラムである。このデーモンは REST API を提供しているのでローカルからだけでなくネットワーク経由でのコンテナの操作が可能である。主要機能は以下の通りである。

- 非特権コンテナ、リソース制限を用いたセキュアなデザイン
- 直感的なコマンドラインと REST API
- イメージベースのコンテナ構築
- ライブマイグレーション

特に、Docker Hub にあるイメージを利用可能になるということがアナウンスされている点が期待できる。

3.2.5 LXD インストール

Ubuntu ユーザは PPA を使って以下の通りインストール可能である。なお、他のディストリビューションのユーザは、最新のリリースの tarball か git リポジトリから直接 LXD をダウンロードしてビルドできる。

```
$ sudo add-apt-repository ppa:ubuntu-lxc/lxd-git-master
$ sudo apt-get update && sudo apt-get -y install lxd
```

3.2.6 LXD イメージのインポート

イメージベースなので、ダウンロードする。なお、LXD のコマンドラインは `lxc` というコマンドである。何ともややこしい。コンテナイメージのインポートは `lxc-images` というコマンドを利用する。以下では、Ubuntu14.04 をインポートしている。

```
$ sudo lxd-images import lxc ubuntu trusty amd64 --alias ubuntu
```

以下のコマンドでイメージ一覧を取得できる。

```
$ sudo lxc image list
+-----+-----+-----+-----+-----+-----+
| ALIAS | FINGERPRINT | PUBLIC | DESCRIPTION | ARCH |          UPLOAD DATE          |
+-----+-----+-----+-----+-----+-----+
| ubuntu | 04aac4257341 | no      |              | x86_64 | Jul 15, 2015 at 1:16pm (UTC) |
+-----+-----+-----+-----+-----+-----+
```


3.2.7 LXD コンテナの起動

`lxc launch` コマンドで起動できる。

```
$ sudo lxc launch ubuntu test-container-102
Creating container...done
Starting container...done
error: saving config file for the container failed
```

できなかった。。。 どうやらこのバグにヒットしたらしい。(https://github.com/lxc/lxd/issues/739)
改めて起動するとこのようになる。

```
$ lxc list
+-----+-----+-----+-----+-----+-----+
|      NAME      | STATE |   IPV4   |   IPV6 | EPHEMERAL | SNAPSHOTS |
+-----+-----+-----+-----+-----+-----+
| test-container-103 | RUNNING | 10.0.3.138 |      | NO        | 0          |
+-----+-----+-----+-----+-----+-----+
```

このような感じで LXD を扱えるが、まだまだバグもあり不安定だという印象が強い。もし、興味があれば使ってみてほしい。



第 4 章

あとがき

4.1 こじろー

本体は表紙です。ついでに、中身も流し読みしていただけると嬉しいです。

4.2 まっきー

hogehoge

4.3 だーまり

hogehoge

