

Security of Docker containers

Report for the Computer Security exam at the Politecnico di Torino

Carmine D'Amico (239540)

tutor: Antonio Lioy

July 2018

Contents

1	Introduction	2
2	State of the art	3
2.1	From virtual machines...	3
2.2	...to containers	3
2.3	LXC	4
2.3.1	Kernel namespace	4
2.3.2	Cgroups	5
2.4	Docker	5
2.4.1	History	6
2.4.2	Docker's architecture[7]	6
2.4.3	Docker's objects[8]	7
2.4.4	Security by default	8
3	Hardening Docker	9

1 Introduction

Explain here why the XYZ protocol is important and what was the purpose of the present work.

If you want to reference a web site you can do like this: <http://www.polito.it>.

2 State of the art

In computer science, the term *virtualisation* is referred to the creation of virtual computational resources. These resources, normally supplied as hardware, are instead provided to the user by the operating system through the creation of a new abstraction layer. OSs, storage devices or network resources could all be virtualised. Virtualisation can be obtained at different levels and using different techniques.

Virtual machines have represented for many years the state of the art of virtualisation, being used in both consumer and enterprise contexts. In the last years a new technology, based on *containers*, has started to gain more attention, thanks to its benefits. Docker is an open source container technology, stepped into the limelight thanks to its simple interface, which allows to create and control containers.

2.1 From virtual machines...

With the term virtual machines it is often intended an **hypervisor-based virtualisation**, that is a type of virtualisation that acts at hardware level. Virtual machines (VMs) are established on top of the host operating system, providing applications with their dependencies, but also an entire guest OS and a separate kernel. One or more virtual machines can be run on the same machine. Hypervisors are distinguished in two different types, the one that works directly on top of the host's hardware (**bare metal hypervisor**) and the one that is on top of the host's OS (**hosted hypervisor**).

Bare metal hypervisor provides better performances, not having the overhead of the extra layer of the host's operating system. It manages directly hardware and the guest's operating system. On the contrary hosted hypervisor can be managed in an easier way, running as a normal computer program on the user's operating system.

As said before, the hypervisor needs to run on the user's computer, which is defined as *host machine*, while each virtual machines is called *guest machine*. It is important to remember this terminology, because it will be used also in the following, referring to containers.

2.2 ...to containers

Container-based virtualisation represents another approach to virtualisation, mainly spread in the last years. Compared to hypervisor-based virtualisation it results more lightweight, using the host's kernel to run multiple virtual environments. It virtualises at operating system level (it is also known as **OS-level virtualisation**), allowing other applications to run without installing their own kernel on the host. Containers look like separated process that just share host's kernel.

Resources are provided by the host's OS, together with the container engine. A container engine is the technology in charge of create and manage containers. **Docker** represents one of the most important and most used container engine. A computer program running into a container can only see the resources allocated to that particular container. On the same host there could be more than one container, each one with its personal set of dedicated resources. Although it could be possible to run more than one computer program inside the same container, it is always suggested to run only one program per container, in order to separate areas of concern. It is better to connect multiple containers using user-defined networks and shared volumes.

There are many representatives of this container-based virtualisation approach, like *Linux-VServer*, *OpenVZ* and *Linux Container (LXC)*. This last will be better described in the next section, allowing to better understand the behaviour of containers.

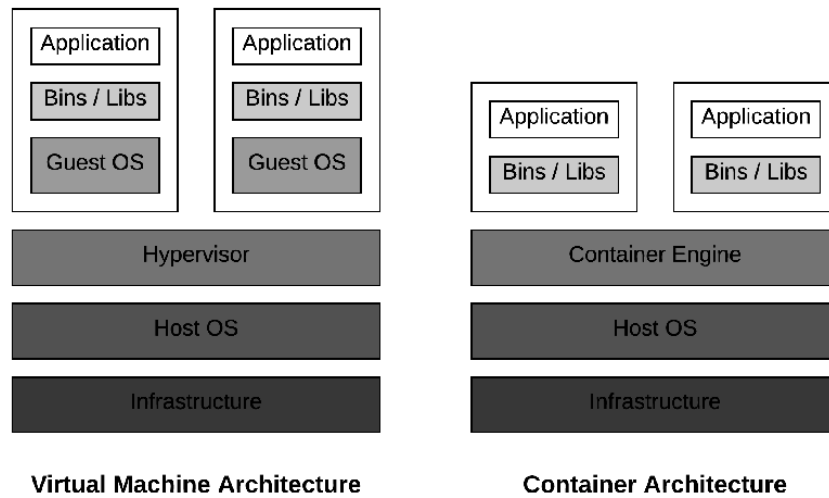


Figure 1: Architecture differences between virtual machine and container

2.3 LXC

Linux Containers, better known as LXC, are an OS-level virtualisation technique created around 2008. They allow to run multiple isolated Linux instances (the *containers*), on top of a single host *LXC host*, which shares its Linux kernel. Each container *sees* its own CPU, memory, network interface, I/O, ecc...

The isolation among containers is obtained thanks to some Linux kernel's tools: namespace and cgroups. In the following subsections these two components will be analysed, both for their importance for containerisation in general and for the fact that they are also the basic components in Docker.

2.3.1 Kernel namespace

Namespaces allow to create isolated environments, in which each process that belongs to that particular environment can see global host's resources as personal isolated resources. In other words they allow to create pool of processes that think to be the only ones of the system. In this way groups of processes, that are part of different namespaces, can see different set of resources. Namespaces work by assigning to different resources the same name in different namespaces. In the Linux kernel six different type of environments are implemented:

- PID Namespace
- IPC Namespace
- UTS Namespace
- Network Namespace

- User Namespace
- Mount Namespace

2.3.2 Cgroups

Cgroups are a kernel tool used to manage processes' resources. They gather, track and limit processes' usage of resources. It is possible to create and manage *cgroups* using high level code, assigning PID to a specific *cgroup*. They represent the fundamental tool to obtain resource isolation, playing an important role also for the CPU and I/O's scheduling. The resources that can be limited by Cgroups are[1]:

- **memory** - this subsystem sets limits on memory use by tasks in a cgroup and generates automatic reports on memory resources used by those tasks.
- **CPU** - this subsystem uses the scheduler to provide cgroup tasks access to the CPU.
- **CPUacct** - this subsystem generates automatic reports on CPU resources used by tasks in a cgroup.
- **CPUsets** - this subsystem assigns individual CPUs (on a multicore system) and memory nodes to tasks in a cgroup.
- **blkio** - this subsystem sets limits on input/output access to and from block devices such as physical drives (disk, solid state, or USB).
- **net_cls** - this subsystem tags network packets with a class identifier (classid) that allows the Linux traffic controller (tc) to identify packets originating from a particular cgroup task.
- **net_prio** - this subsystem provides a way to dynamically set the priority of network traffic per network interface.
- **ns** - the namespace subsystem.
- **devices** - this subsystem allows or denies access to devices by tasks in a cgroup.
- **freezer** - this subsystem suspends or resumes tasks in a cgroup.
- **perf_event** - this subsystem identifies cgroup membership of tasks and can be used for performance analysis.

2.4 Docker

As today, *Docker* represents the most used computer program for operating-system-level virtualisation (containerisation). It is developed by **Docker, Inc**[2] and it was introduced during the 2013's PyCon. During its presentation,

Docker was announced as the future of Linux Containers[3], indeed in its first releases it reiterated many concepts from them, such as *Namespaces* and *Cgroups*, but providing a simpler user experience and a complete ecosystem to create and manage containers.

Docker's success is mainly addressable to its portability and lightweight nature, that allow to create high density environments. It is the ideal software in scenarios where continuous

integration and continuous delivery (CI/CD) are required, allowing developers to not only build their code, but also test their code in any environment type and as often as possible to catch bugs early in the applications development life cycle [5].

2.4.1 History

Docker was born as an inside project within *dotCloud*, a platform-as-a-service, later renamed to **Docker**. Solomon Hykes[4] was the leader of the project, that was at first developed with other *dotCloud*'s engineers, like Andrea Luzzardi and Francois-Xavier Bourlet. The project went public, as said before, during 2013's PyCon and it was released as open source software during the same year. Always during 2013, *Docker* distanced itself from LXC (exactly with version 0.9), replacing it with a new execution environment, **libcontainer**.

Docker represented a turning point in the IT industry, as it can be proved by looking at its adoption. The following is a list of the milestones achieved by the program from Wikipedia[6]:

- On September 19, 2013, Red Hat and Docker announced a collaboration around Fedora, Red Hat Enterprise Linux, and OpenShift.
- In November 2014 Docker container services were announced for the Amazon Elastic Compute Cloud (EC2).
- On November 10, 2014, Docker announced a partnership with Stratoscale.
- On December 4, 2014, IBM announced a strategic partnership with Docker that enables Docker to integrate more closely with the IBM Cloud.
- On June 22, 2015, Docker and several other companies announced that they are working on a new vendor and operating-system-independent standard for software containers.
- As of October 24, 2015, the project had over 25,600 GitHub stars (making it the 20th most-starred GitHub project), over 6,800 forks, and nearly 1,100 contributors.
- A May 2016 analysis showed the following organisations as main contributors to Docker: The Docker team, Cisco, Google, Huawei, IBM, Microsoft, and Red Hat.
- On October 4, 2016, Solomon Hykes announced InfraKit as a new self-healing container infrastructure effort for Docker container environments.
- A January 2017 analysis of LinkedIn profile mentions showed Docker presence grew by 160% in 2016.[30] The software has been downloaded more than 13 billion times as of 2017.

2.4.2 Docker's architecture[7]

Docker follows a client-server architecture, where three main components can be distinguished (Fig. 2):

- The server, which is a daemon process (called **dockerd**) running on the host's machine. It is in charge of create and manage Docker's objects.
- A set of interfaces conformed to REST architectural style, that enable programs to communicate with the server, sending instructions.

- A command line interface (CLI) client, that allows the user to interact with Docker using the REST API through their terminal.

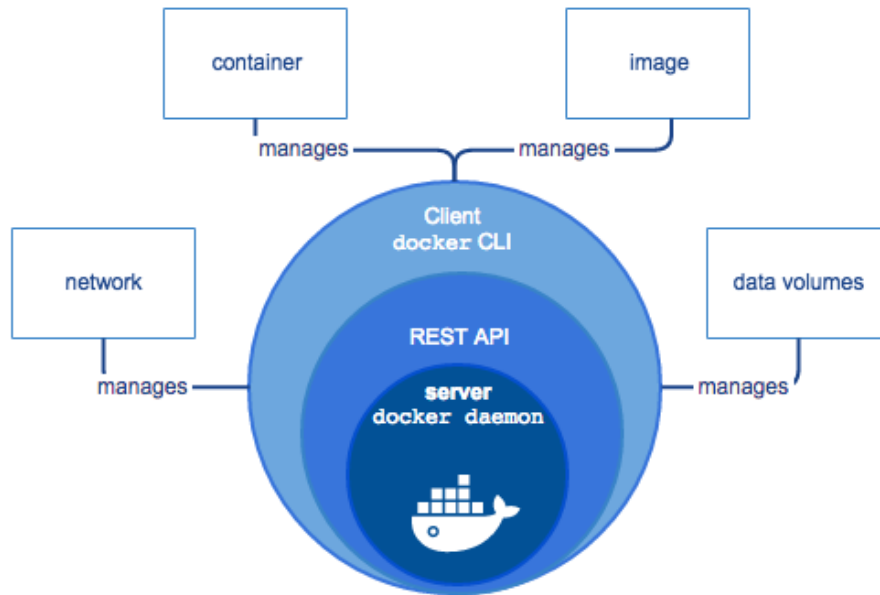


Figure 2: Docker's architecture.

2.4.3 Docker's objects[8]

Docker's workflow includes the interaction with many special purpose objects, created and managed by **containerd**:

- **IMAGES** are the basic components involved in the creation of a Docker's container, containing all the instructions that the daemon has to follow. An image can be created from scratch or it can be based on already existing images (for example on the image of *nginx*), where other needed components are installed. A *Dockerfile* is a special file, that follows a very simple syntax, which includes all the steps that must be followed in order to create an image. Each instruction represents a layer in the image. When a new layer is inserted (modifying the *Dockerfile*) and a container is rebuilt, only the new layer is rebuilt, speeding up the deployment's process.

Docker images are stored inside registries, that can be private or public. The two most famous public registries are *Docker Cloud* and *Docker Hub*, the latter is the default one visited by Docker for searching images.

- **CONTAINERS** represent the the runnable instances of images. The relationship between images and containers could be compared to the relationship between classes and objects in an object-oriented programming language. A container can be connected to the network or to a storage and it can be defined by its image or by the configurations indicated starting it.

When a container is started Docker searches locally for all the needed images (downloading them from online public registries if necessary), then a read/write filesystem is allocated,

where the container can create or modifies file or directories. By default a container can be connected to the external network using the host's connection. When a container is stopped any changes to its state that are not stored in persistent storage disappear.

- **SERVICES** are supported from version 1.12 of Docker. They allow to run containers across different daemons. A series of daemons connected between them compose a *swarm* and they all communicate between them using Docker's REST API. A user can defines *swarm*'s configurations, like the number of *service*'s replica available. A *service* is seen from the external as a single application.

2.4.4 Security by default

3 Hardening Docker

References

- [1] https://access.redhat.com/documentation/en-us/red_hat_enterprise_linux/6/html/resource_management_guide/ch01
- [2] <https://www.docker.com/>
- [3] <https://www.youtube.com/watch?v=wW9CAH9nSLs>
- [4] https://en.wikipedia.org/wiki/Solomon_Hykes
- [5] <https://www.docker.com/use-cases/cicd>
- [6] [https://en.wikipedia.org/wiki/Docker_\(software\)#History](https://en.wikipedia.org/wiki/Docker_(software)#History)
- [7] <https://docs.docker.com/engine/docker-overview/#docker-architecture>
- [8] <https://docs.docker.com/engine/docker-overview/#docker-objects>