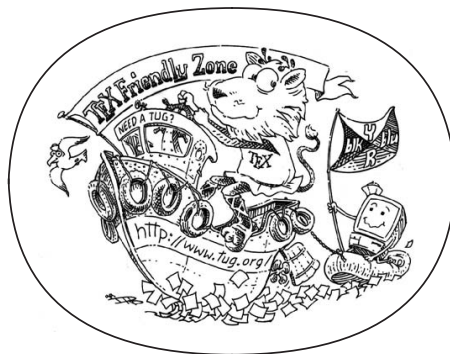


ANDRÉ MIEDE & IVO PLETIKOSIĆ
A CLASSIC THESIS STYLE

A CLASSIC THESIS STYLE

ANDRÉ MIEDE & IVO PLETIKOSIĆ



An Homage to The Elements of Typographic Style

June 2018 – classicthesis v4.6

POLITECNICO DI MILANO
V Facoltà di Ingegneria
Corso di laurea in Ingegneria delle Telecomunicazioni
Dipartimento di Elettronica e Informazione

TITOLO DELLA TESI

Relatore: Prof. Nome COGNOME

Correlatore: Ing. Nome COGNOME

Tesi di Laurea di:
Simone Mosciatti Matr. 878758

Anno Accademico 2017-2018

Ohana means family.
Family means nobody gets left behind, or forgotten.
— Lilo & Stitch

Dedicated to the loving memory of Rudolf Miede.
1939–2005

ABSTRACT

Short summary of the contents in English... a great guide by Kent Beck how to write good abstracts can be found here:

<https://plg.uwaterloo.ca/~migod/research/beck00PSLA.html>

ZUSAMMENFASSUNG

Kurze Zusammenfassung des Inhaltes in deutscher Sprache...

*We have seen that computer programming is an art,
because it applies accumulated knowledge to the world,
because it requires skill and ingenuity, and especially
because it produces objects of beauty.*

— knuth:1974 [knuth:1974]

ACKNOWLEDGMENTS

Put your acknowledgments here.

Many thanks to everybody who already sent me a postcard!

Regarding the typography and other help, many thanks go to Marco Kuhlmann, Philipp Lehman, Lothar Schlesier, Jim Young, Lorenzo Pantieri and Enrico Gregorio¹, Jörg Sommer, Joachim Köstler, Daniel Gottschlag, Denis Aydin, Paride Legovini, Steffen Prochnow, Nicolas Repp, Hinrich Harms, Roland Winkler, Jörg Weber, Henri Menke, Claus Lahiri, Clemens Niederberger, Stefano Bragaglia, Jörn Hees, Scott Lowe, Dave Howcroft, José M. Alcaide, David Carlisle, Ulrike Fischer, Hugues de Lassus, Csaba Hajdu, Dave Howcroft, and the whole L^AT_EX-community for support, ideas and some great software.

Regarding L_YX: The L_YX port was initially done by *Nicholas Mariette* in March 2009 and continued by *Ivo Pletikosić* in 2011. Thank you very much for your work and for the contributions to the original style.

¹ Members of GuIT (Gruppo Italiano Utilizzatori di T_EX e L^AT_EX)

CONTENTS

1	INTRODUCTION	1
2	BACKGROUND	3
2.1	WLCG	3
2.2	CVMFS	5
2.2.1	CVMFS High Level Overview	5
2.2.2	CVMFS Details	6
2.3	Containers	6
2.3.1	Docker and the cvmfs/graphdriver plugin	7
2.3.2	Singularity	9
2.4	Problem Definition	9
3	STATE OF THE ART	13
3.1	Distributed File System	13
3.2	Provising of Machines	14
3.2.1	Single o Local Machines and Small Clusters	14
3.2.2	Local cluster	15
3.2.3	WLCG environment	16
3.2.4	Managing run time dependencies on the WLCG	17
3.3	Working lazily	18
4	METHODOLOGY	21
4.1	High Level overview of the proposed file system	21
4.2	Singularity Images	23
4.3	Docker Thin Images	25
4.4	Keeping track of the work already done	27
4.5	Closing remarks	29
5	IMPLEMENTATION	31
5.1	CVMFS Write Interface	31
5.2	Singularity Ingestion	31
5.3	Docker Ingestion	32
5.3.1	CVMFS Ingestion of Tarball	33
5.3.2	Docker Ingestion Algorithm	33
5.4	Garbage Collection of Images	34
5.5	Administrator Interface	34
5.6	The <i>repository-manager</i> command line tool	36
6	RESULTS	37
6.1	Container Startup Time	37
6.2	Space Requirement	38
6.3	File in the sub-catalogs	39
6.4	Complexity	40
7	FUTURE WORKS	43
A	APPENDIX	45
A.1	Code to get the numbers	45

BIBLIOGRAPHY	49
--------------	----

LIST OF FIGURES

Figure 2.1	Schematic representation of the WLCG	4
Figure 2.2	Decision proces for running docker thin images	10
Figure 4.1	High level visualization of the proposed file-system	22
Figure 4.2	Visualization of the Filesystem structure, the arrows indicate symbolic links	23
Figure 4.3	Visualization of the "super directories" in the ".flat" subdirectory	24
Figure 4.4	Complete visualization of the .flat directory	28
Figure 4.5	Structure of the .metadata/ directory	28
Figure 6.1	Startup times of the containers served using the proposed architecture and served using the standard Docker Registries.	39
Figure 6.2	Number of files inside the catalogs in the .layer directory.	41
Figure 6.3	Number of files inside the catalogs in the .flat directory.	41

LIST OF TABLES

Table 5.1	Available placeholders and their application to the image https://registry.hub.docker.com/library/centos:centos6	36
Table 6.1	Benchmark of startup time of a containers, the first number is the average while the second is the standard deviation. The units are in hundredths of seconds. $n = 100$	38
Table 6.2	Apparent size in GB of the two folder .layers and .flat	39
Table 6.3	Result of the cyclomatic complexity analysis, only function with complexity greater or equal to 3 are shown	42

LISTINGS

Listing 2.1	Example of a <i>recipe</i> file of a Docker <i>thin image</i>	8
Listing 4.1	Algorithm to add an image reference to the layer metadata	26
Listing 4.2	Algorithm to remove an image from the file-system	27
Listing 5.1	Algorithm to unpack a Docker image with Singularity and ingest it into CVMFS	32
Listing 5.2	Example of a small <i>wish-list</i>	35
Listing 6.1	Storage require to store the whole repository	38
Listing A.1	Script used to capture the startup time of singularity with image hostes in CVMFS using CVMFS cache	45
Listing A.2	Script used to capture the startup time of singularity with image hostes in CVMFS without cache	45
Listing A.3	Script used to capture the startup time of docker thin-images using both CVMFS and Docker cache	45
Listing A.4	Script used to capture the startup time of docker thin-images without Docker nor CVMFS cache	45
Listing A.5	Script used to capture the startup time of Docker standard images without Docker cache	46
Listing A.6	Script used to capture the startup time of Singularity running Docker standard images without Singularity cache	46
Listing A.7	Script used to capture the startup time of Singularity running images unpacked on the local file-system, hence with cache	47
Listing A.8	Script used to capture the startup time of docker thin-images without cache	47

ACRONYMS

INTRODUCTION

CERN, European Organization for Nuclear Research, (French: Conseil européen pour la recherche nucléaire) is an European research organization that operates the largest particle physics laboratory in the world.

Its mission is to:

- Provide a unique range of particle accelerator facilities that enable research at the forefront of human knowledge
- Perform world-class research in fundamental physics
- Unite people from all over the world to push the frontiers of sciences and technology, for the benefits of all.

It host the instrument of the 4 biggest physic collaboration:

- ALICE
- ATLAS
- CMS
- LHCb

CERN host also a plethora of smaller physical collaborations that benefits from the instruments, know how, network effect and services availables.

Between the services offered to its users the computing service is one of the most interesting. Indeed CERN host and manage one of the biggest computing data center used for public research.

An issues that affect the operations inside the data center is the provisioning of software on the computing servers.

A specialization of the same problem is the provisioning of containers images.

The general problem of software provisioning is been solved by the use of CernVM-FileSystem [4], a read only file-system that provides a scalable, reliable and low maintenance software distribution system.

This thesis will explore the problem of creating a suitable read-only file-system structure to provision containers images on computing nodes. We will provide a general read-only file-system structure and we will implement the proposed methodology on top of CVMFS.

This thesis is composed by several parts: The background will provide the necessary information on the CERN computing architecture (WLCG), then we will explore CVMFS and why it is a good fit for the

CERN computing architecture, the last part of the background will cover the integration between CVMFS and containers technologies. The state of the art will explore what alternatives are available for software distribution in general and for distribution of images. We will then define the problem that this thesis is trying to solve and few metrics of interest in our specific case. The methodology part will explain the details of the solution we propose for this specific problem. The implementation chapter will focus on how the proposed methodology is been put in practise. We will evaluate the result of the proposed methodology and implementation on the result part following the metrics that were previously proposed. Finally we will propose future work and enhancement to the implementation

BACKGROUND

In this chapter we are going to introduce the concept and the technologies that made this work possible. We will start introducing the Worldwide LHC Computing Grid (WLCG), which is the *collaboration* that provide the computing power necessary to the CERN mission. The dimension of the WLCG and its specific workload required and allowed a specific software distribution system, CernVM-FileSystem which is used to provision the machines on the WLCG.

Then we will introduce the concept of containers, a different way to solve the software distribution problem that is widely adopted in the industry. In particular, we will focus on Docker containers and the Docker images format. We will then explore Singularity, a container runtime capable of running containers images stored as simple folder in the filesystem. Finally we will explore the Docker `cvmfs/graphdriver`, a Docker plugin that allow to run Docker images whose content is available on a read-only file-system without the need of downloading or accessing the Docker standard images.

WLCG

The Worldwide LHC Computing Grid is an global collaboration of more than 170 data centers in 42 countries. The mission of the WLCG is to provide the computing resources to store, distribute and analyze the data generated by the operations of the LHC [6],[8].

The organization of the WLCG follows a hierarchical model, where each level of the hierarchy is called “Tier”. The most central Tier is the Tier-0, which is hosted by CERN in the Geneva Area and in Budapest. There are 13 Tier-1 data centers with enough storage and computing capabilities to support the Grids operation around the clock. Tiers-1 are geographically distributed: 8 of them are in Europe, 3 in North America and the rest is in Asia. Finally, the Tier-2 data centers do not have strict requirements and are generally operated by research centers and universities [8]. A schematic representation of the architecture of the WLCG is provided on Figure 2.1.

The work of the WLCG is mostly divided in two big classes: analysis of the data from the LHC detectors and Monte Carlo simulations. The WLCG assume and support a batch computing paradigm. Analysis and simulations are split in smaller jobs that are distributed to different computing node that can work in parallel [8], [3].

Before to start each job it is necessary to install the software on the server. Unfortunately the amount of software potentially needed in

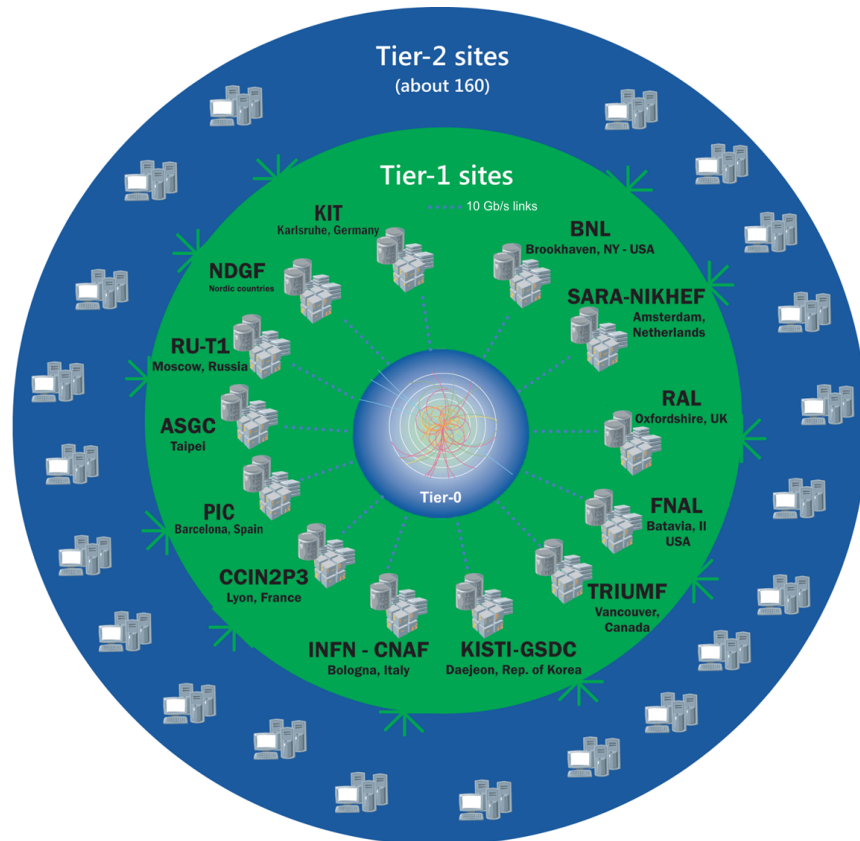


Figure 2.1: Schematic representation of the WLCG

each computing node and the velocity at which the software is updated can make the installation challenging. Moreover, simpler installation techniques that rely on packages managers are not applicable since they would put the centralized package managers themselves under too much load. Several solution have been proposed and used in the past, eventually it settled for the use of CernVM-FileSystem (CVMFS) [4].

CVMFS

This section will explore CernVM-FileSystem, we will start with an high level overview of CVMFS, then we will explore what happens when a file is requested.

CVMFS High Level Overview

CernVM-FileSystem [4] provides a scalable, reliable and low-maintenance software distribution system. It is implemented as a read-only POSIX file-system in user space exploiting FUSE (File-system in USErspace) [26] and standard web server technologies such as Apache or NGNIX.

Each running instance of CVMFS provides a read-only file-system that is denominated *repository*. At CERN different collaborations maintains different repositories, but all of them can be mounted from all the computing node in the WLCG. CVMFS is engineered to support repository of size on the order of the Terabyte with billions of files.

To save storage space files are addressed by their content (Content Addressable Storage), hence duplicated files will be stored only once.

In order to distribute software to geographically distant data centers and keep a low latency, CVMFS allows to cache content in different machines. This allow to host a cache server in each Tier of the WLCG. The use of caches fits perfectly with the Tiers models of the WLCG presented above. The Tier-0 host the main repository (Stratum-0), and the Tier-1 host the first level of cache (Stratum-1) and so on.

The content of the files are served using the HTTP protocol by a standard web server. The files are lazily downloaded only on the machine that need them and only when necessary.

In order to locate and request files from CVMFS the clients download the catalog, a simple SQLite database which describes a subtree of the whole file-system. The catalog contains all the metadata of files and directories, including owner, group, permission, and size. Moreover the catalog contains also the URL where to download the files.

A root catalog is available in a know path, and, if the file-system grows too large, the root catalog links to other sub-catalogs. The use of sub-catalogs allows to keep each catalog small improving the query time.

CVMFS Details

CVMFS is implemented using the Client-Server architecture. The server is responsible to manage the content of the repository and to expose it via HTTP API. The client is installed in the host machine and is responsible to expose the content of the repository to the users and it is implemented as a FUSE daemon which implements all the system calls necessary for a read-only file-system.

When a CVMFS file-system is mounted, it starts by reading a configuration file which describes each repository. The client then downloads a simple text file which points to the catalog of the repository. Once the catalog is downloaded the client has all the information necessary to start responding to the system calls performed by the user.

As an example, when the user requires a `stat` system call against a file, the client reads from the catalog all the information about the file like, size, permission, mode, etc.. and replies with them. Instead when the user requires to read from a file, the client first downloads the file from the server, stores it into a local cache and passes through each read operation to the local copy.

This approach allows to download only the file really required, since all the other system calls can be served by just reading from the catalog. However this implies that the reading latency from a file depends on the network latency. This strategy works very well if the reading latency of a file is not a major concern and if reading from the catalog is fast. However, if the catalog grows too big, then the queries become too slow.

To overcome this limitation the sub-catalogs were introduced. A sub-catalog is exactly like the normal catalog, but while the catalog refers to the whole file-system tree, a sub-catalog refer to a smaller sub tree of the file-system. In order to avoid confusion, we will refer to the root-catalog as the catalog that includes the root of the file-system and to sub-catalog to all the other catalogs in the file-system. The root-catalog, of course, embeds several sub-catalogs, and each sub-catalog can, recursively, embed another sub-catalog.

When is required to read information about a file, the client starts by looking into the root catalog and then it follows the sub-catalogs structure until it finds the required file.

CONTAINERS

While CERN solved its problems of software distribution with CVMFS the industry opted for a different approach: containers. Containers are a standard units of software that package up code and all its dependencies so that computer applications run quickly and reliably from one computing environment to another [19].¹

In order to standardize containers and make the technology interoperable in 2015 the *Open Container Initiative (OCI)* was founded [11]. The OCI defined a standard format to use to pack containers into *images* [oci-image-spec], this format have been adopted by Docker and that is used in the *Docker Images*. A Docker image is an immutable set of tar files, where each tar file is called layer. Prior to run the container, each layer gets mounted one on top of the other to re-create the original environment where to run the application [9]. The content of each image is codified in a json file, the manifest, which provides the unique name of the image itself and which refers to each layer that compose the image by their unique identifiers. The unique identifier of both images and layers is the result of the function hash256 of their content [10].

Docker images are distributed through Docker registries, simple HTTP servers that given the unique identifiers of a layer provide the layer itself, similarly, given the identifier of an image the registries provides its manifest [16].

Docker allow to associate a human readable identifier to each image, this name is composed by a namespace, which identifies the user or organization that created the image, a name, which identifies the image itself and a tag, which identifies the version of the image. These names are not immutable and are meant to be used just by humans to recognize and use the images [18]. The repository where the image is hosted, its namespace, its name and its tag create a hierarchical structure between the several images that is easy to navigate for humans.

Docker and the cvmfs/graphdriver plugin

Docker is a thin CLI layer on top of a Linux daemon, dockerd which is responsible to obtain the Docker images from the registries, mount the layers, manage the runtime of the image and allow communication of a Docker runtime with the host system [17].

Docker layers can be shared between different images, so the Docker daemon download them once and store them in the local machine for future use [14]. On average less than 7% of the content inside a Docker image is used during the runtime of the image [13]. The layers are distributed as a single tar file [9], hence a naive distribution of layers with CVMFS would not provide any advantage. The whole tar files would need to be downloaded erasing the advantages of using CVMFS.

A solution to this problem was introduced with the cvmfs/graphdriver plugin and the concept of "thin-images" [12]. In this work we are not interested in the internal of the cvmfs/graphdriver plugin, it will be sufficient to know that the Docker daemon can be enhance with the

use of plugins [15] and that the plugin allows us to run what are called "thin-images".

Thin-images are Docker images that are created starting from standard (*fat*) Docker images. The content of a *thin-image* is only a simple json file whose content is the list of layers needed by the original *fat-image* and where to find them in the host file-system, we call this file the *recipe* an example of *recipe* file is show on Listing 2.1. *Thin-images* can be executed only if the Docker daemon have enabled the cvmfs/graphdriver plugin, indeed the plugin is able to read the content of the thin-image, mount the original layers inside the container, and finally execute the application [12].

Listing 2.1: Example of a *recipe* file of a Docker *thin image*

```
{
  "version": "1.0",
  "origin": "https://registry.hub.docker.com/library/python:
    latest",
  "layers": [
    {
      "digest": "bc9ab7...0257b7",
      "url": "cvmfs://thin.osg.cern.ch/.layers/bc/bc9ab7...0257b
        7/layerfs"
    },
    {
      "digest": "193a63...110d7f",
      "url": "cvmfs://thin.osg.cern.ch/.layers/19/193a63...110d7f
        /layerfs"
    },
    {
      "digest": "e5c3f8...818580",
      "url": "cvmfs://thin.osg.cern.ch/.layers/e5/e5c3f
        8...818580/layerfs"
    },
    {
      "digest": "b61a0d...b615a9",
      "url": "cvmfs://thin.osg.cern.ch/.layers/b6/b61a0d...b615a
        9/layerfs"
    }
  ]
}
```

If the files necessary to run the *thin-images* are distributed by CVMFS, it is possible to run docker containers without downloading any unnecessary content.

Moreover the cvmfs/graphdriver plugin is able to run standard (*fat*) Docker images, hence conventional Docker images are supported seamlessly [12].

On figure 2.2 we show how the plugin decides how to mount the container file system. If it detects that the image is a standard *fat*

Docker image it run it as a standard image, if it detect that is a *thin* image it mounts the layers from the CVMFS read-only file system.

The downside of this approach is that a "thin-image" can be downloaded, store in the local storage and successfully run using the content provide by CVMFS. On a second moment the images can be update, uploaded newly on the registry and the old content in CVMFS is deleted to host the new content of the image. If now the user try to run the same images, it will not download it again from the registry but it will use the image already cached, referring to content not available anymore on CVMFS, hence it will fail to run.

Singularity

Singularity [29] is another container runtime, it provides its own image format but is capable to run standard Docker images [28] as well. Moreover it is capable of running containers, also Docker containers, directly from a directory containing the unpacked container file-system itself [27].

Given the Singularity capability to run Docker containers directly from a directory containing the unpacked Docker image file-system its integration with CVMFS is simpler that the one with Docker itself. Indeed it is sufficient to host the directory containing the unpacked file-system of the Docker image in CVMFS.

PROBLEM DEFINITION

We have introduced how CERN has overcome the challenges of software distribution using CVMFS. Unfortunately this is not enough since run-time dependencies between software components can break applications that instead are running perfectly fine in a different environment.

Since containers pack up all their runtime dependencies in a standard environment they are a suitable solution for this problem. However, containers are distributed as big tar files which is against the design, working principles and performance characteristics of CVMFS.

Both CVMFS and containers aim to solve the same challenge of server provisioning, however they made different trade-offs. CVMFS opted for an efficient distribution of the content, making difficult and inconvenient to pack all the runtime dependencies of a particular application. Containers, instead, opted to make each application self-contained so that the application can run reliably on different computer environments, however losing efficiency in the distribution of the content.

Efficient content distribution and efficient dependency management should not be contrasting goals. For this reason, in this thesis work

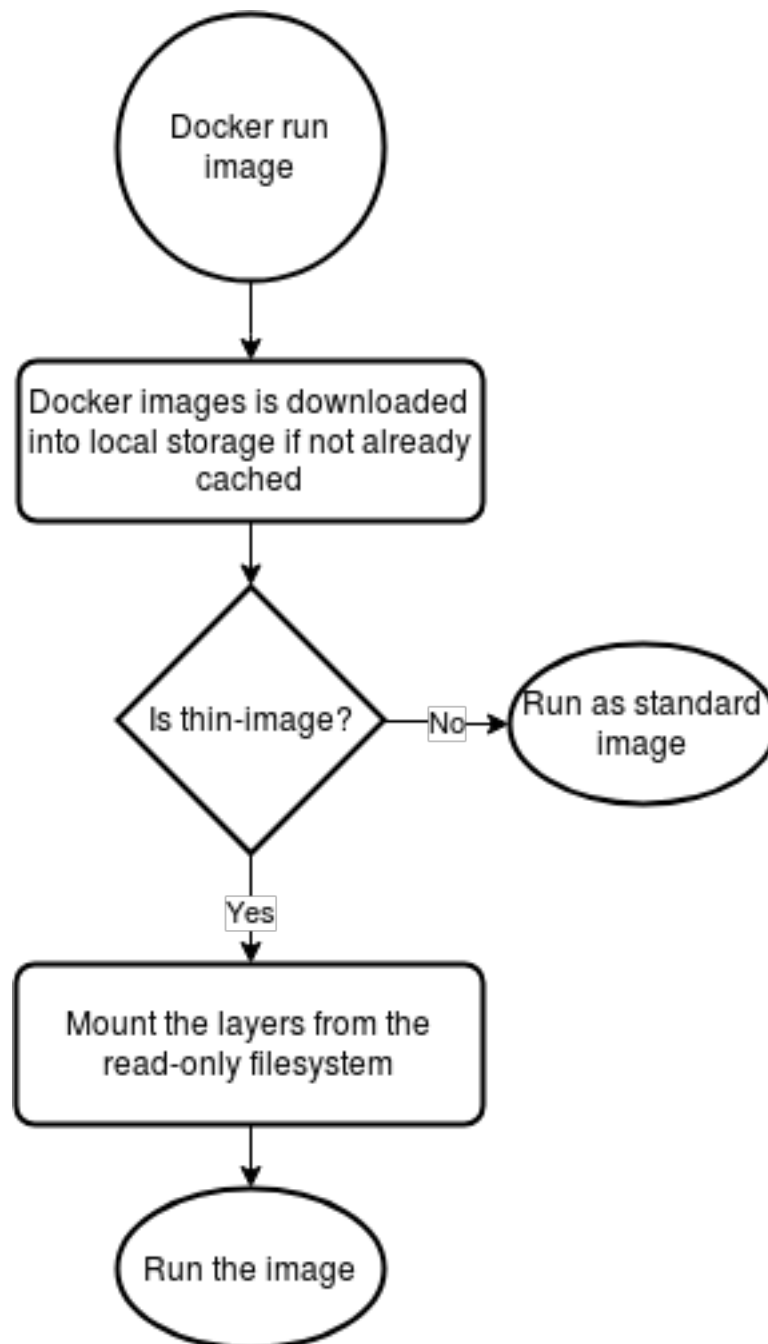


Figure 2.2: Decision proces for running docker thin images

we will tackle the problem of how to efficiently distribute software in HPC clusters while efficiently preserving all its runtime dependencies.

STATE OF THE ART

In this chapter we provide the state of the art for the technologies that are of interest for this work. In Section 3.1 we will discuss the state of the art in distributed file system where we will see the major difference between CVMFS and other distribute file system. In Section 3.2 we will explore how provisioning and installation are done on computing clusters of various size, starting from the smallest up to the globally distributed one. Finally in Section 3.3 we will explore early achievement of working lazily when provision machines.

DISTRIBUTED FILE SYSTEM

Several distributed file system have been proposed and are available. One of the first implementation of a distributed file system is the Andrew File System (AFS) [22]. It provides a read-write file system, however it guarantee only weak consistency, hence read and write operations are directed only to locally cached copy of the file. The changed are propagated to the server when a file is eventually closed. AFS inform the client if a different user modify a cached file using callbacks, but those callbacks are discarded after any communication error between the client and the server, including network failures or timeouts. The need to keep a constant connection open between the client and the server makes the client/server ration quite small, 200. Hence for every 200 clients a new server machine is needed.

Another commercial distributed file system is BeeGFS it has a different approach than AFS, since it is composed by three different components, the clients, the metadata server and the storage servers. The clients coordinate and informs the metadata servers which in turns coordinate the storage servers that, finally, are accessed by the clients. Also this architecture of distributed file system relies on metadata components that bound the number of clients imposing an overhead.

Other approach to store huge amount of data produced Amazon S3 (Simple Storage Service). S3 does not provide a classical POSIX interface but an eventually consistent REST interface. Objects are created immediately but deletions will eventually happen in the future moreover it has no knowledge of permissions outside of a simple access control list, symbolic links or even folders. Due to this shortcomings S3 is not directly usable to distribute and instal software one HPC clusters, however it can be used as a simple layer to access the file leaving all the metadata handling to CVMFS. As mentioned in 2.2 the files provide by CVMFS are accessed using the HTTP protocol, they

don't have to be in a web server managed by CVMFS but they can be in any accessible location.

Up to our knowledge no distributed file system handle metadata like CVMFS does, indeed all the metadata queries to CVMFS are handle by the client themselves reading the catalogs files, hence the metadata operations are not a limitation anymore since they scale out with the number of client themselves. Of course this is possible since CVMFS provides to the clients only a read interface and the write are all serialized in the CVMFS server component.

PROVISING OF MACHINES

In this section we are going to explore how it is possible to provision and install software on cluster of machines. The proposed methodologies will grow in complexity as the size of the cluster grows. In parallel we will see how the problem shifts from the handling of files to the handling of metadata.

In Section 3.2.1 we will describe how installation and provision of machine is done in humble environemnt with very few machines. The need of automation will arise on bigger cluster, hence in Section 3.2.2 we will describe how automatic tools help in provisioning local clusters. In Section 3.2.3 we will explore the solution that were applied to the WLCG before CVMFS and their weakness. Finally in Section 3.2.4 we will explore proposed solution to the problem of managing run time dependencies in a globally distributed cluster like the WLCG.

Single o Local Machines and Small Clusters

Several system have been proposed to manage installation and provision of machines.

The simpler approach is of course the manual approach were the source code is downloaded on the local machine, compiled and finally installed. While this works fine for the expert user, it introduces a lot of difficulties and friction for users that are not experts. The compilation part is error prone and it may requires dependencies, so the user must be ready to doing a lot of research and even reading makefiles if necessary. Moreover, if the software has run-time dependencies, it is necessary for the user to install them as well, and so recursively.

Of course there is no need to re-compile the software every time, moreover compilation is an expensive and quite technical step that not all users are confident in doing. Indeed the package managers solves this kind of problem, the software is compiled once by expert on the software itself and distributed along with a script that install it in the most appropriate location in the system and with a series of dependencies. This is a major shift in the interface since now the user simply installs a package without caring of all its dependencies

or without even needing a compiler. It is the package manager that will eventually install run time dependencies, directly or indirectly.

The distribution of packages poses already a challenge, if enough users are interested in the packages hosted in a single package manager, it would be necessary to load balance the request between several mirrors site, indeed this is what happens in big linux distribution that serves a lot of users. The package manager itself is mirrored to different data center each one serving a part of the traffic.

Local cluster

As long as the installation of software is manual and done in relatively small clusters package managers are suitable for hosting all the software and provision all the servers.

When the cluster start to grow, manual installation is not a solution anymore, indeed specialized tools are available to manage the state of a fleet of machines.

Ansible [20], [2] is a tool that helps in managing the installation of software in a fleet of machine. It is based on two simple concepts, the *inventory* and the *playbook*. The inventory associate each machine that is managed to a set of tags. The playbook is a set of statuses each associated with a tag. When provisioning a fleet of machines, Ansible tries to adapt the status of each server in a way that it respect the status declared in the playbook. As an example in the playbook we may declare that the servers tagged as web-server need to have the apache service on, while the servers tagged as build-node needs to have the gcc package installed. In the first case Ansible will try to start the apache service while in the second case it will try to install the gcc package.

A possible issue with Ansible is about the use of those servers, users may login into the server, change their configuration to fix an urgent problem and then don't updates the Ansible playbooks. In the long term this causes a drift on configuration and settings that makes quite difficult to manage the fleet of services.

A solution to the drift in configuration is the use of a different tool, *Puppet* [25], [2] which works on a similar principle of Ansible but it constantly monitor the status of each machine reverting any changes.

If different users need to access a shared pool of computing resource other solution consist in the use of *virtual machines*. Each physical machine get partitioned into a set of different, isolated virtual machines where each users can install and configure its software. Moreover the image of the virtual machine itself may comes with software already installed and configured.

Providing software already installed in a machine is a major conceptual shift where we move from the user having the responsibility to install all the needed software to a model where the software was

already in the system itself. Virtual machines while providing the stronger level of isolation possible in a shared environment are slow to set up and also slower than native performance.

An evolution of virtual machines are the containers: containers provide a smaller level of isolation than virtual machine but are much quicker to set up and have close to native performance. Moreover containers are simple to build and their popularity is raising between developer as well. Indeed it is possible to set up a container that encapsulates all the status necessary to run a specific computation, this allow to move the container from a machine to another without worrying anymore of run-time dependencies.

All these solutions works well in local clusters, however all of them requires to move all the necessary software on a local machine, also components that are not needed. Starting a fleet-wide installation using tools like *Ansible* or *Puppet* means that the package managers needs to serve at roughly the same time all the packages necessary by all the machines, which can quickly sum up to a lot of bandwidth. The use of virtual machines or containers is not better in this cases since it still requires to move a lot of content. In big cluster is common to provide local package managers and local container registries.

Providing the content to several different machines is a problem that can be solved by careful use of cache and optimization. Diving deeper into the containers worlds the OCI standard defines that each layer is identified by a single unique digest. This helps since it provides an unique identifier to cache and, more importantly, by bundling together a lot of different files in a single tar files it sidestep the need of managing metadata for all the files in a container. Still this comes at the cost of needing to download a big file even if its content is not all necessary. Similar consideration can be done for the provisioning of Virtual Machines.

WLCG environment

With the respect of provisioning the WLCG and its modus operandi is an harsher environment that local cluster. Indeed it operates in a batch computation mode, where jobs are scheduled and started all together in a wide fleets of machines, possibly in geographical different data centers. A provisioning structures that relies on package managers or simple containers will witness an huge peak in the bandwidth request that will slow down the provisioning itself. Moreover scientific software is usually quite large, also in the order of 20G and the provisioning of relatively few machines is sufficient to cause a big spike in the bandwidth request. Indeed several solution have been proposed and tested to address this issue.

The *Andrew File System* [22] (AFS) is a distributed read-write file system optimized for home directories on globally distributed work-

stations. AFS provides a global name space (/afs) partitioned into *cell* and *volumes*. Scientific collaborations can host the authoritative copies of their experiment software on public volumens (cell /afs/cern.ch). Using AFS, scientific software need to be installed only once in the server in order to be used by several clients. Since AFS needs to keep the connection state for all its clients it has a low client/server ration with is limited to about 200:1.

Another proposed mechanism to manage installation on the WLCG are the *grid installation jobs*, described in [5]. Jobs that pre-install software releases on worker node. In order to avoid to install the software on every single node WLCG Tiers provides a *share software areas*, a NFS volume mounted to every worker node that can be used to provide the software itself. While the *share software areas* are a necessity, they also introduce a single point of failure and high chance to overload the NFS with meta-data operations.

An evolution of the *grid installation jobs* is the *ALICE Software Installation System* described in [5]. The software releases are still distributed as packages but those packages are physically distributed to the worker nodes using the BitTorrent protocol. This improves the installation times using all the bandwidth not only from the package manager through the worker nodes but also between the worker node themselves. Unfortunately this architecture required to still transfer all the packages to all the nodes which is a waste of bandwidth, introduce sources of errors and it is time consuming.

Another evolution of the *grid installation jobs* is *GROW-FS* [7] which pioneered the idea of meta-data handling on the worker nodes and not in a central node. *GROW-FS* provides as well a shared areas for installation using a simple web server. *GROW-FS* uses a single catalog for the entire directory tree and changes to the file system structure requires to completely reconstruct the catalog as well as re-mounting the file system on the worker nodes.

Managing run time dependencies on the WLCG

Other than the distribution of content, CERN faced the problem of managing run-time dependencies. A possible solution to this problem would be to statically link all the dependencies, but this would generate extremely big binaries. Another solution would be to carefully managing all the software installed on the machine, including all the recursive dependencies. However this clash with the WLCG model where jobs can migrate from a data center to another along with their dependencies. Moreover sometimes even a very careful installation it is not sufficient since is possible that two application that could run on the same machine have clashing dependencies.

The problem of how to manage runtime dependency is been solved outside CERN in several different way.

On small scale it is possible to use simple package managers that automatically install the dependencies, however they suffer of few different problem. First and foremost in case of dependencies clashing the package manager simply refuse to install the required package, which is not an acceptable solution. Then at the scale of the WLCG the package manager quickly become a bottleneck with respect to the bandwidth necessary to distribute the content but also with respect to the size of the internal databases.

Another possible solution is the use of containers. Packaging all the runtime dependencies along the application code in an immutable file system allow containers to exactly reproduce the same environment and condition to ensure smooth operations. The use of containers however is not convenient in very large clusters like the WLCG, indeed the content necessary to run a container is distributed as few large tar files and this makes the distribution itself inefficient since a lot of content not useful for the application itself is downloaded anyway.

Several system have been propose to decrease the downloading time of containers images, however all of them work *ahead of runtime* trying to optimize how fast the overall system is capable of delivering the whole image into the host. FID (Faster Image Distribution) [21] is a P2P Docker images distribution system that is able to accelerate the speed of distributing Docker images by taking full advantage of the bandwidth of not only the Docker Registry but also of the other nodes in the cluster, while decreasing the downloading time. However, FID still require to download all the image before to run the computation. Another work by Anwar et al. [1] characterize the workload of large-scale registries in order to derive design implications for more optimized registries, still they are working to optimize the time necessary to serve and download the whole image from the registry.

While the aim of this work is still to decrease the start up time of uncached containers we took a different road. The previous works focus on optimizing the delivery time of the whole content image, we decided to focus on minimizing the amount of content that the client needs to download which of course leads to shorter start up time for the containers and also to save bandwidth.

WORKING LAZILY

The experience with the installation mechanisms tested at CERN suggests that it is necessary to work lazily in order to provide software and content to a wide fleet of machines and this is exactly what CVMFS does for us. However we need another step in order to solve also the problem of run time dependencies.

Few other system have been proposed to address the problem using a similar architecture. The *cvmfs/graphdriver* [12] allows us to

run Docker images starting from a *recipe* file that list the layers to mount and their location in the file system. Unfortunately a way to provide the *thin images* was not provided, moreover it was not provide a way to structure the file system itself. Slacker [13] tackle the same problem but it uses a very different implementation. Their work is based on NFS and flatten layers. All the layers of an image get flattened into a single layer, and such layers stored as a single file into a NFS shared between the Docker Registries and the Docker Client. While *cvmfs/graphdriver* shares a set of layers as Docker images, Slacker share a simple reference to the file stored in the NFS. Slacker is able to reach very interesting performance thanks to the implementation of the NFS they use that relies on network disks. We lack the details about the NFS used by Slacker but we believe that their performance may be strongly influenced by the size of the cluster.

Similarly to Slacker in another work Nicolae et al. [24] propose another lazy structure this time to run Virtual Machines. They trap the IO to the Virtual Machine image using FUSE and redirect those call into a shared read-write file system backed by local disk on each machine. Similarly to Slacker we believe that this approach works well on a small set of machines but it won't scale well on a big cluster. Indeed their experimental set up was limited to 120 machines. As mentioned above for the *grid installation jobs* the use of NFS brings the risk of overloading the NFS itself with metadata operations.

In this work we aim to finally bridge the two worlds of containers and distributed read-only file system like CVMFS. Instead of working with blocks like Slacker and Nicolae et al. our primitives will be files. Similarly we are not going to provide a read-write interface to the worker node since this will definitely impact the scalability of the solution given the necessities to implements locks. The client will only be able to read the files.

In this chapter we will introduce a read-only file-system structure for running Docker images using both Singularity and the docker thin-images plugin. We focus on Singularity because it is a widely deployed system in HPC and on Docker thin-images because it allows the user to leverage the Docker infrastructure to run images whose content is distributed with a read only file-system.

Section 4.1 will provide a high level overview of the proposed methodology explaining why we decided to focus our design on CVMFS, Singularity and Docker with the *cvmfs/graphdriver* plugin.

Section 4.2 will analyze the file-system structure to host the Docker images to run with Singularity. We will explain how we created a hierarchical structure similar to the one of the docker registries while keeping the repository maintainable and without putting too much pressure in the sub-catalog system of CVMFS.

Section 4.3 will analyze how we stored the content used by the thin-images Docker plugin. Structural similarities between Docker images and Docker layers will drive us to adopt a very similar solution to avoid stressing the sub-catalog system. The sharing of layers between different docker images will allow us to avoid repeating work, but it will introduce difficulties during the deletion of the image itself that we will solve using a reference count system.

Finally Section 4.4 will show how we kept track of the images already present in the file-system.

HIGH LEVEL OVERVIEW OF THE PROPOSED FILE SYSTEM

As mention the aim of this work is to provide an efficient content distribution system to run containerized application on big distributed cluster.

The aim of this work is how to efficiently distribute software in HPC clusters while efficiently preserving its runtime dependencies. To tackle this problem we will resort on CVMFS to efficiently distribute the content of the containers together with containers technologies to preserving the runtime dependencies in such a way that only the strictly necessary files will be downloaded into the worker nodes of the clusters.

The content distribution part will be managed by CVMFS which provides us with the primitives necessary to efficiently distribute the content. Moreover, to the least of our knowledge, CVMFS is the only stable and tested system that allows to distribute lazily and efficiently

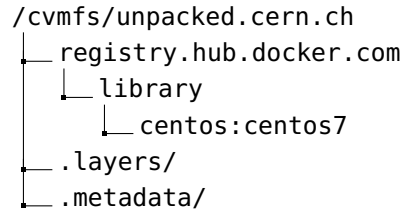


Figure 4.1: High level visualization of the proposed file-system

Terabytes of data in a distributed data center like the WLCG. Starting from the data distributed by CVMFS, there are two container run-times that we can use to actually run the containers. The first one is Singularity, introduced in Section 2.3.2 that allows to run containers whose content is already unpacked in a simpler folder. The second one is Docker with the *cvmfs/graphdriver* plugin, which allows Docker to run *thin images* made of a simple *recipe* json file which describes where are in the file system the necessary layers that are needed to be mounted before to start the container itself.

The proposed structure is absolutely generic and even if not implemented with CVMFS it would allow the above container run-times to work with minimal modifications ¹. However some choices of the directory structures have been taken starting from the design of CVMFS and may not be necessary on a generic read-only file system.

The proposed file-system provides two main structure. A non-hidden set of directories that host the unpacked images to run with Singularity and a hidden folder to host the layers used by the *cvmfs/graphdriver* Docker plugin. Along with these structures another hidden directory will contains metadata information about the Docker images already in the repository.

The directories that host the unpacked Docker images have the same hierarchical structure of the Docker structure mentioned in 2.3, hence the first directory is the name of the registry that host the image, it follow the namespace which refers to the user of the organization responsible for the image and the last level is the image name along with the tag.

The directory that host the layers of the Docker images can conceptually be a flat structure simply containing the layers each identified by its digest. However a simple flat structure will put too much pressure in the catalog system so we aggregated layers that share the same digest prefix and create a sub-catalog for each aggregation.

The hidden metadata folder will follow the same hierarchical structure of Docker images, allowing to quickly locate the metadata information of an image given just the name of the image itself.

¹ The implementation of *cvmfs/graphdriver* relies on a specific *cvmfs* url in the *recipe* format which is a simple implementation details that can be easily generalized.

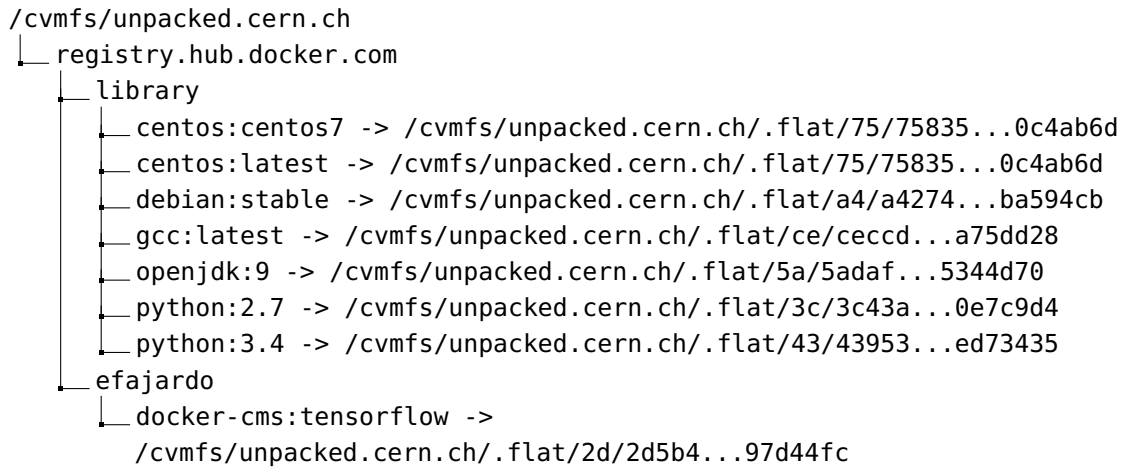


Figure 4.2: Visualization of the Filesystem structure, the arrows indicate symbolic links

SINGULARITY IMAGES

To run Docker images using Singularity is sufficient to start the singularity executable providing as input the directory where the image is been unpacked. In this section we are going to show how we structure the file-system in a way that allow users to easily discover and run unpacked docker images using Singularity while keeping the file-system easy to maintain.

As mentioned in 2.3 docker images have a hierarchical structure. The first level of the hierarchy is the docker registry where the image is hosted. The most common registries in our case are the official docker hub ² and the CERN internal registry ³.

The second level in the hierarchical structure is the namespace of the docker image. If the image is one of the official docker images it will be the standard namespace: “library”. In all the other cases, the namespace will be the same as the original docker image. For example for the images belonging to the ATLAS collaboration we use the namespace `atlas`.

The last level is the name of the image itself together with the tag of such image, separated by a colon (:). We decided to avoid yet another level containing just the tags. Indeed there are relatively few tags for each image and adding another level of indirection would have made it harder to explore the file-system. Moreover, we decided to use the colon because it is the same character used in the docker registries between the images and the tag and it is immediately recognizable by users.

While this structure is user friendly, it makes the maintenance of the repository complex.

² registry.hub.docker.com

³ gitlab-registry.cern.ch

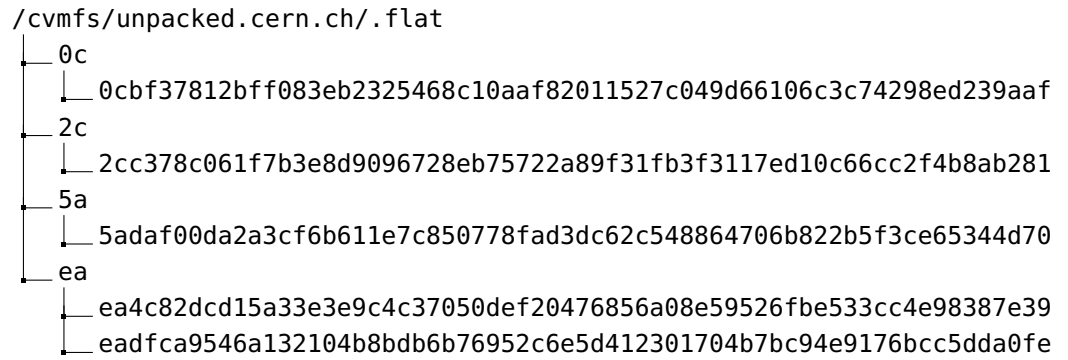


Figure 4.3: Visualization of the "super directories" in the ".flat" subdirectory

The tags used in each image are not immutable, hence, without continuous maintenance, it may happen that the images stored inside the file-system are not up to date making difficult for the user to know what version of the software is being run. However with the described structure, it would be extremely complex to detect if an image is up to date or if it needs further updates.

To work around this issues we exploited the fact that each image is uniquely identified by its digest. Indeed we decided to store the real content of the images in an hidden folder that embed the digest itself while preserving the structure presented above using symbolic links.

We show the directory structure of the file-system on figure 4.2.

The folder that contains the real content of a Singularity images are all below the standard subdirectory `.flat/`. The name `.flat/` was chosen to make it clear that only flattened file systems are stored in there.

Embedding the digest in the name of the folder allows to immediately find the location of an image, which is useful when an image become obsolete and need to be deleted from the file-system.

From a theoretical point of view it would be sufficient to store the whole content of the Singularity images in the folder `.flat/$image_digest`. However, from a practical point of view this would create too much content in a single folder putting too much pressure in the CVMFS sub-catalog system.

To overcome this issue we decided to create a fixed number of "super-directories" where we placed the unpacked folder of the images. To easily locate each unpacked folder in the super-directories we decided to call each super-directory as the prefix of the digest of the images it is containing. Since the digest is an hexadecimal string this approach provides us with $16 \times 16 = 256$ fixed super-directories inside the `.flat/` directory, each of which will contain only the content of the images whose digest start with those 2 specific bytes.

On figure 4.3 we can see that "0c", "2c", ..., "ea" are all "super-directories" and each one contains only the file-systems that start with

“oc”, “2c”, ..., “ea” respectively. Note the case of “ea” that contains file-systems of multiple images whose digest start with “ea”.

However, to relieve pressure from the catalog system is not sufficient to simply aggregate the images into “super-directories”, we also need to create a sub-catalog for each “super-directory.” Moreover, since each image can contains itself a lot of files we decide to create a sub-catalog also for each unpacked image.

Another positive side-effect of the use of symbolic links is that symbolic links manipulation is defined as atomic in the POSIX standard.

The use of “super-directories” is necessary for limits in the implementation of CVMFS and they are not necessary on an abstract read-only file-system.

DOCKER THIN IMAGES

While for running docker images using Singularity it is sufficient to have the image unpacked in a simple directory running docker containers using the thin-images plugin requires a more complex set up. As explained in 2.3.1 the recipe of the docker thin-image contains the path of the directories where each layer of the original docker image is hosted, those directories will be mounted by the docker plugin.

All the docker layers are stored under a common subdirectory of the file-system, the `.layers/` directory.

Since the sub-tree of the file-system used by the Docker thin-images is used only by the Docker plugin we don’t need to create a human-friendly structure like we did for the Singularity sub-tree.

Like docker images also the docker layers are identified by an unique digest, and similarly to the docker images, store all of them in a single directory will put too much pressure in the CVMFS sub-catalog system, hence we follow the exact same model used for storing the unpacked images also for the layers creating 216 super-directories.

A big advantage of the use of layers over flat images, is that layers can be shared by multiple images.

The sharing of layers allow us to avoid re-doing work that is already been done, in particular if a later is already in the file-system it will not be added again. On the other hand it makes more complex removing an image since it is necessary to remove each layer that compose the image, but some layer may be shared between images.

Removing layers has the important implication that once the layer is removed every thin image that relies on it won’t work anymore. However those thin-images could be stored on the client side where we don’t have any access.

To do not disrupt the user workflow while keeping the repository to a manageable size we consider several option:

1. Never remove layers

Listing 4.1: Algorithm to add an image reference to the layer metadata

```

Function AddReferenceToImage
    Pass In: LayerReference, ImageReference
    ReferenceFile := FindReferenceFile(LayerReference)
    if ReferenceFile exist
        References := LoadReferenceFromFile ReferenceFile
        Add ImageReference to References
        Overwrite References to ReferenceFile
    else
        References = ImageReferences
        Write References to ReferenceFile
    endif
EndFunction

```

2. Remove layers as soon as possible
3. Provide a grace period before finally removing the layer

The option to never remove layers is impractical since the size of the file-system will grow unbounded.

Remove layers as soon as possible is not desirable, even running computation could be broken by this policy and the users have no way to deal with this possibility but retrying the whole computation.

The last option is the most sensible and better suited for our use case, and so it is the one that we implement, this gave users the possibility to:

1. Complete their computation
2. Update the local images in order to always run stable containers

In order to know which layer to delete from the file-system we store a reference that map each layer to the images that use the layer itself.

These references are stored as metadata in a simple .json file. We store one of these reference file for each layer in the file-system.

Anytime a new image is added to the file-system we update the several references files, adding for each layer in the image, a reference to the image itself.

When we decide to remove an image, for any layer we check that it is used only by the image we want to remove, if this is the case, we remove the layer, if it is not the case we just remove the reference of the image.

In order to store both the metadata information about the layers (in particular the "reference" file mentioned above) and the actual file-system of the layer an additional directory structure is used. Below the directory called as the digest of the layer there are two more directories:

1. layerfs/ directory that actually store the content of the layer

Listing 4.2: Algorithm to remove an image from the file-system

```

Function RemoveLayer
    Pass In: LayerReference, ImageReference
    ReferenceFile := FindReferenceFile(LayerReference)
    References := LoadReferenceFromFile ReferenceFile
    Remove ImageReference from References
    if size References == 0
        Remove Layer
    else
        Overwrite References to ReferenceFile
    endif
EndFunction

```

2. `.metadata/` directory that stores the references to the image in a simple JSON encoded file, “`origin.json`”

Of course, the recipe of the thin images is not concerned at all with the content of the `.metadata/` directory. Hence the recipe files points directly to the `layerfs/` directory.

The complete structure for storing docker images is the one showed in [4.4](#)

KEEPING TRACK OF THE WORK ALREADY DONE

To avoid to perform duplicated work is necessary to keep track of which image is already been added to the file-system. The same information may be used by the users to know exactly what images is hosted in the file-system.

In order to know which image is already been converted we need to uniquely identify each image, as already mention, using the combination of image name and tag is not enough, since the tag are mutable. Hence we rely on the digest of the image.

The information about each image is stored into another top-level hidden directory, `.metadata/`.

Inside the `.metadata/` folder we have others directories, one for each hosted image. Inside those directories there is a single file, `manifest.json` that store the manifest of the image itself.

As already mentioned in [2.3](#) the manifest contains the digest of the image itself. Comparing the manifest stored in the file-system with the manifest downloaded from the docker registries is possible to understand if the image should be updated or nor.

The structure of the `.metadata/` folder is show in figure [4.5](#).

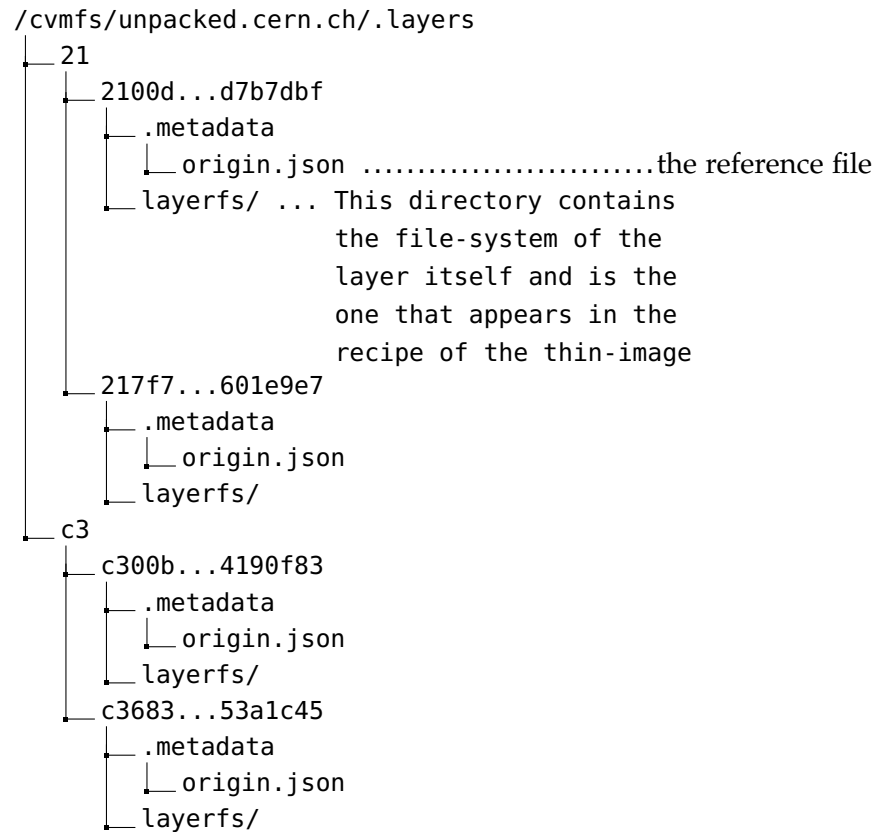


Figure 4.4: Complete visualization of the .flat directory

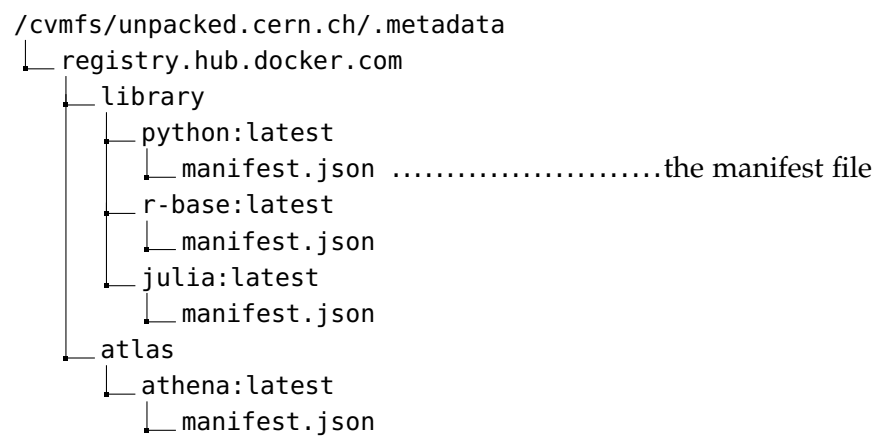


Figure 4.5: Structure of the .metadata/ directory

CLOSING REMARKS

In this chapter we have introduced a file-system structure suitable to host docker images that can be run using both Singularity and the docker thin-images plugin.

We started by storing the unpacked images used by Singularity in a hierarchical structure that recall the one of the Docker registries to enhance the discoverability of the images itself. This approach however would have make difficult to maintain the repository since we would not know the version of each image unpacked. We overcome this issue storing the real unpacked images in a hidden directory that embed the digest of the image itself and using symbolic links to preserve the hierarchical structure. Too many unpacked images stored under the same directories however would have put too much pressure on the catalog system of CVMFS, hence we adopted the use of super-directories.

On the second part we analyzed how to store the layers used by the docker thin-image plugin. A single layer can be used by multiple images this allowed us to avoid repeating work but at the same time it makes more complex removing an image from the file-system. We decided to keep a reference count to know when is safe to actually delete a layer. The same issue of too many files under the same sub-catalog arose also for storing the layers, we used the same approach used for the unpacked images based on the use of super-directories.

The last section explored how we keep track of exactly which image is already store in the file-system storing the image catalog in a hidden subdirectory.

IMPLEMENTATION

In this chapter we are going to explore how the ingestion inside CVMFS is been implemented. At first we will describe the write interface of CVMFS, then we are going to talk how we unpack Docker images to use them with Singularity. We will move on to describe how we ingest Docker *fat images* while transforming them into Docker *thin-images* suitable to be used with the *cvmfs/graphdriver* plugin, moreover we will explore the custom modification done to CVMFS in order to accommodate the needs of layer ingestion. We will then explore how all the bookkeeping of images and layers is manages. Then we are going to understand how the images are removed from the file system. Finally we will show the interface given to the administrator.

All the work presented in this chapter is implemented in the *repository-manager*, a command line utility written in the Go(lang) language [23].

CVMFS WRITE INTERFACE

As mentioned in Section 2.2.2 CVMFS is implemented with a Client-Server architecture: while the client is strictly read-only, the server does provide a write interface. What we write in the server is what the client is then able to read. The write interface of CVMFS is a transactional interface, hence it is possible to open a transaction, modify the file system and then either publish the modification or abort the transaction. Those actions are carried out respectively by the command ‘*cvmfs_server transaction*’, ‘*cvmfs_server publish*’ and ‘*cvmfs_server abort*’. When a transaction is open the server file system is a standard writable linux file system, hence it is possible to modify it using the standard POSIX API, Linux commands or even graphical file explorers. Moreover it is possible to test locally the new file system without actually commit the changes. Of course these actions are not available to the clients that have access only to a read-only interface.

Finally it is important to keep in mind that only a single transaction can be opened at any given time since CVMFS will refuse to open a second transaction.

SINGULARITY INGESTION

The step to ingest a singularity images are pretty straightforward. Initially the image is downloaded from the remote registry and stored in a temporary area. The download is carried out by Singularity itself, in

Listing 5.1: Algorithm to unpack a Docker image with Singularity and ingest it into CVMFS

```

Function UnpackAndIngestDockerImage
    Pass In: DockerImageName

    Digest := RetrieveDigestFromImageName DockerImageName
    TemporaryDirectory := CreateTempDirectory
    UnpackDockerImageWithSingularity TemporaryDirectory

    StartCVMFSTransaction
        FlatDir := CreateFlatDirectory Digest
        MoveFrom TemporaryDirectory Into FlatDir
    CommitCVMFSTransaction

    HumanReadableName := GetDirectoryFromDockerImageName
        DockerImageName

    StartCVMFSTransaction
        CreateDirectory HumanReadableName
        CreatSymlinkFrom HumanReadableName To FlatDir
    CommitCVMFSTransaction

EndFunction

```

order to minimize the possibilities of inconsistencies or of errors. Once the download completed successfully and the unpacked container file-system is in the local file-system we start the real ingestion phase.

The first step of the ingestion is to open a transaction in CVMFS. Once we open the transaction we copy the temporary directory into the CVMFS filesystem under the “.flat/” directory. Then we commit the first transaction. A second transaction takes care of creating the symbolic link as described in 4.2.

This few steps are sufficient to make the Singularity images available through CVMFS. The pseudocode 5.1 show the details of the algorithm.

DOCKER INGESTION

Converting Docker *fat images* into Docker *thin images* is a more complex task than simply make the unpacked image available to use with Singularity. We will explore all the details of this process in this section.

In the first part we will introduce how CVMFS is able to directly ingest tar file which is the format used to distribute Docker images as mention in Section 2.3. Then we will explore how starting from the Docker manifest we add the layers to the CVMFS file system, then we show how we create the *thin image* itself.

Moreover we also upload the Docker *thin image* to a Docker registry.

CVMFS Ingestion of Tarball

As mention in 2.3 Docker layers are distributed as tar files. In order to support the use case of ingesting Docker layer we decided to add a new command to CVMFS ‘cvmfs_server ingest’. The ingest command takes as input a tar file and a directory inside the CVMFS file system and extracts all the files and directories in the tar file into the directory provided as input. This command implicitly opens and commits a CVMFS transaction, hence it is possible to have only a single concurrent ingestion.

Docker Ingestion Algorithm

The first step of the algorithm is to retrieve the manifest of the Docker images from the Docker registry, as soon as we have the manifest the algorithm checks if the specific image is already been successfully ingested into the file system. This check happens using the metadata stored in the *.metadata* directory as mention in Section 4.4. The check consists in a simple comparison between the digest of the manifest of Docker images just downloaded and the digest of the images already stored in the *.metadata* folder. If the images is already in the file system the algorithm terminates.

The next step is to ingest each layer of the Docker image into the CVMFS file system. As previously we check if the layers already exist in the file system itself. Since the layers are stored under a path that embeds theirs own digest as describe in Section 4.3, checking if a layer is already in the file system consists in simply checking if the folder where we would ingest the layer already exists or not. If the layer already exists we move to the next layer of the image.

The ingestion of a layer follows a similar procedure of the ingestion of an unpacked image, the layer is first downloaded into a temporary directory and then it is ingested using the ‘cvmfs_server ingest’ command. Another option could have been to avoid storing the layer in the temporary directory and simply let the ‘cvmfs_serve ingest’ command read the content of tar file from STDIN, we decided against this approach since a non negligible amount of times the download of the layer fails in the middle wasting all the work already done by the ‘ingest’ command.

If the ingestion of any of the layers fails we stop the whole algorithm and we rely on retries from the administrator in order to have the Docker image served on CVMFS.

After all the layers have been successfully ingested the next step of the algorithm is the creation of the Docker *thin image*. The Docker *thin image* is a standard Docker image which content is a single file

`thin.json` that contains the set of layers to mount before to start the container itself as mentioned in Section 2.3.1. To create this image it is sufficient to encode the location of the layers into a `.json` file and then pack this file into a standard Docker image. The Docker *thin image* is then uploaded into the Docker registry.

The last step of the algorithm is to store the metadata information about the image just ingested in order to avoid to repeat work already done. This is done simply storing the manifest of the docker image in the `.metadata` folder following the schema presented in Section 4.4.

GARBAGE COLLECTION OF IMAGES

We have described how we can add new images to the CVMFS file system, however updating an image is quite common, especially if the images are referred by mutable tag such as “latest” which actually represent the latest version of a particular application.

During the update of an image we avoid to immediately delete the files from the CVMFS repository, as mentioned in Section 4.3 this could cause disruption of service for users.

Instead we keep track of all the images that are not necessary anymore in a specific file, the `remove-schedule.json` file which is stored in the hidden directory `.metadata/` just below the main root of the CVMFS file system. The `remove-schedule.json` files contains a collections of the manifest of all the images we are not interested in anymore.

When it is time to actually delete all the old images we scan the `remove-schedule.json` file and we carry out the actual removal of the images.

The removal of a singularity image is quite simple, indeed, is sufficient to remove the whole directory.

Removing the layers of the docker images is more complex. At first we need to identify all the layers that we need to check. This is simple since this information is stored in the manifest itself which is stored in the `remove-schedule.json` file.

Then, for each layer we obtain the list of images that need the layer itself. From that list we remove the image we are eliminating from the file system. If the list is now empty we proceed to remove also that specific layer from the file system.

ADMINISTRATOR INTERFACE

In order to store images in the CVMFS file system is necessary to know what image to store, in which repository store it and how to call the respective `thin-image`. We decide to call the triplet `<Input image, CVMFS Repository, Output image>` a *wish*.

To express a list of those wishes we opted for a simple YAML file that store a specialization of a generic “wish list”. In the YAML file we specify a list of Input images, only a single CVMFS repository and a single syntactical transformation for the Output images. An example of this specific *wish list* is show on Listing 5.2.

Listing 5.2: Example of a small *wish-list*

```
version: 1
user: smosciat
cvmfs_repo: 'thin.osg.cern.ch'
output_format: '${scheme}://gitlab-registry.cern.ch/smosciat/thin
               -osg/${image}'
input:
  - 'https://registry.hub.docker.com/library/centos:latest'
  - 'https://registry.hub.docker.com/library/centos:centos6'
  - 'https://registry.hub.docker.com/library/centos:centos7'
  - 'https://registry.hub.docker.com/library/debian:latest'
  - 'https://registry.hub.docker.com/library/debian:stable'
  - 'https://registry.hub.docker.com/library/debian:testing'
  - 'https://registry.hub.docker.com/library/debian:unstable'
  - 'https://registry.hub.docker.com/library/ubuntu:latest'
  - 'https://registry.hub.docker.com/library/fedora:latest'
  - 'https://registry.hub.docker.com/library/python:latest'
  - 'https://registry.hub.docker.com/library/python:2.7'
  - 'https://registry.hub.docker.com/library/python:3.4'
  - 'https://registry.hub.docker.com/library/openjdk:latest'
  - 'https://registry.hub.docker.com/library/openjdk:8'
  - 'https://registry.hub.docker.com/library/openjdk:9'
  - 'https://registry.hub.docker.com/library/gcc:latest'
  - 'https://registry.hub.docker.com/library/r-base:latest'
  - 'https://registry.hub.docker.com/continuumio/anaconda:
    latest'
  - 'https://registry.hub.docker.com/bbockelm/cms:rhel6'
  - 'https://registry.hub.docker.com/bbockelm/cms:rhel7'
  - 'https://registry.hub.docker.com/efajardo/docker-cms:
    tensorflow'
  - 'https://registry.hub.docker.com/lincolnbryant/atlas-wn:
    latest'
```

The syntactical transformation depends on the input image and is applied to obtain the final name of the output image. It simply replace the placeholder \$(PLACEHOLDER) on the Output Image with respective item of the Input Image. We show a reference table on Table 5.1.

The use of a simple YAML file to express the desired content of the file system brings several benefits. Since the wish list can be hosted on a version-control system like Github or Gitlab can be used to host and keep track of the several version of the *wish-list*. Moreover it enables a Pull-Request based approach to change the wish list itself. Users who wish a new image to be added to the repository can simply make a pull request adding their image to the wish list. The administrator

\$(scheme)	https
\$(registry)	registry.hub.docker.com
\$(repository)	library/centos
\$(reference)	6
\$(image)	library/centos:6

Table 5.1: Available placeholders and their application to the image
https://registry.hub.docker.com/library/centos:centos6

of the system act as a gatekeeper, inspect the image that has been required to be added and decides if it is the case to add the image or not.

THE *repository-manager* COMMAND LINE TOOL

All the work presented in this chapter has been implemented in the *repository-manager* [23] command line tool. The tool provide to the administrator of the repository just two main command. The *convert* command and the *garbage-collect* command.

The *convert* command takes as input a *wish list* as show in Section 5.2 and follow all the procedures to add the unpacked Docker images to be used with Singularity to the repository as well as creating the several Docker *thin images* and push them into the Docker registries. Moreover, whenever it needs to delete an image it adds it to the *remove-schedule.json* file.

The actual removing of the images is carried out by the *garbage-collect* command which reads the *remove-schedule.json* files and carried out the removal as described in Section 5.4.

Finally another command of the utility is worth to mention, is the *loop* command, which simply executes the *convert* command in an infinite loop, each time reading again the *wish list* file. This allows the administrator to run the conversion in the background and set up a periodic job that updates the *wish list* file downloading it from a version control system.

RESULTS

In this section we will explore the result of this works, we mentioned 4 metrics that we were considering for this work, first we focus on the start up times of containers, then we analyze the amount of space required in the CVMFS repository and we will compare it with the sum of the size of the containers in the system. We will then analyze the files in each CVMFS catalog and finally we will analyze the complexity of managing the file system itself.

All the measurement were done against a CVMFS repository containing images necessary for standard HEP work.

The measurement of a container start up time is greatly influenced by the caches layers in all cases, either if we serve the content with CVMFS - which use an internal cache - or if we run standard containers where the content may already be cached in the standard file system. Hence we will measure the start up time with and without caches with the content served using both CVMFS and the standard content distribution mechanism of both Singularity and Docker. In each setting we will start the containers 100 times. All the measurement are made inside the CERN data center where we assume a stable and reliable connection.

To quantify the amount of space required we measured the amount of space that the repository uses, comparing this figure with the amount of data that the repository provides. Moreover we will report the total size of all the images store in the repository itself.

The amount of files in each CVMFS catalog is a simple measurement, since the amount of catalogs is rather big we will synthesize this measurement.

Finally to measure the complexity of managing the file system, we measured the cyclomatic complexity of the software that we use to manage it.

CONTAINER STARTUP TIME

This Section explores the startup time of the containers hosted in the CVMFS file-system.

We will show the average time to start a container, start the python shell inside the container and execute the *quit()* command inside the shell. We decide for this simulation because it requires to have access to the python interpreter in the container, which is a executable of not negligible size, 3.4M. Moreover we choose the python images because is a widely know image with a good balance between developing

		Cache		No Cache	
		Avg	STD	Avg	STD
Thin-Image on CVMFS	Singularity	10.31	4.38	29.13	4.02
	Docker	96.74	5.40	180.02	35.13
Naive	Singularity	7.38	2.07	1884.76	366.84
	Docker	95.26	5.29	1279.21	168.99

Table 6.1: Benchmark of startup time of a containers, the first number is the average while the second is the standard deviation. The units are in hundredths of seconds. $n = 100$

work and production work. The size of the whole image is of roughly 923M.

We propose 8 scenarios that we synthesize in the Table 6.1 and on the graph of Figure 6.1. The *Thin-Image on CVMFS* columns represent the startup time of containers which content is hosted in the CVMFS architecture discussed. The *Naive* columns represent the start up time of containers using the default distribution system based on Docker Registries.

We ran the benchmark 100 times, and we collect the startup time with the unix utility `time`. The code that we used to obtain the number is available on the listing on Section A.1

We can see that the use of CVMFS is a huge help when the cache is not available, moreover its overhead is almost negligible when the cache is present. Is interesting to know that the containers hosted in CVMFS without cache starts in an amount of time whose order of magnitude is comparable to the containers hosted naively but using cache (4 times slower in the case of Singularity and 2 times slower in the case of Docker.)

SPACE REQUIREMENT

By default CVMFS stores all its content in `/srv`, hence the most reliable way to obtain the size of the repository is to analyze the storage used under `/srv`. The storage required for the whole repository is of 27G calculate with the `df` unix utility as show in 6.1.

Listing 6.1: Storage require to store the whole repository

```
>> df -h /srv
Filesystem      Size  Used Avail Use% Mounted on
/dev/vdb         197G   27G  161G  15% /srv # <-- I believe is too
               much, most likely an error
```

We than obtain the size of the two hidden folder `.layers` and `.flat` using the `'cvmfs_server list-catalogs'` command which provides us

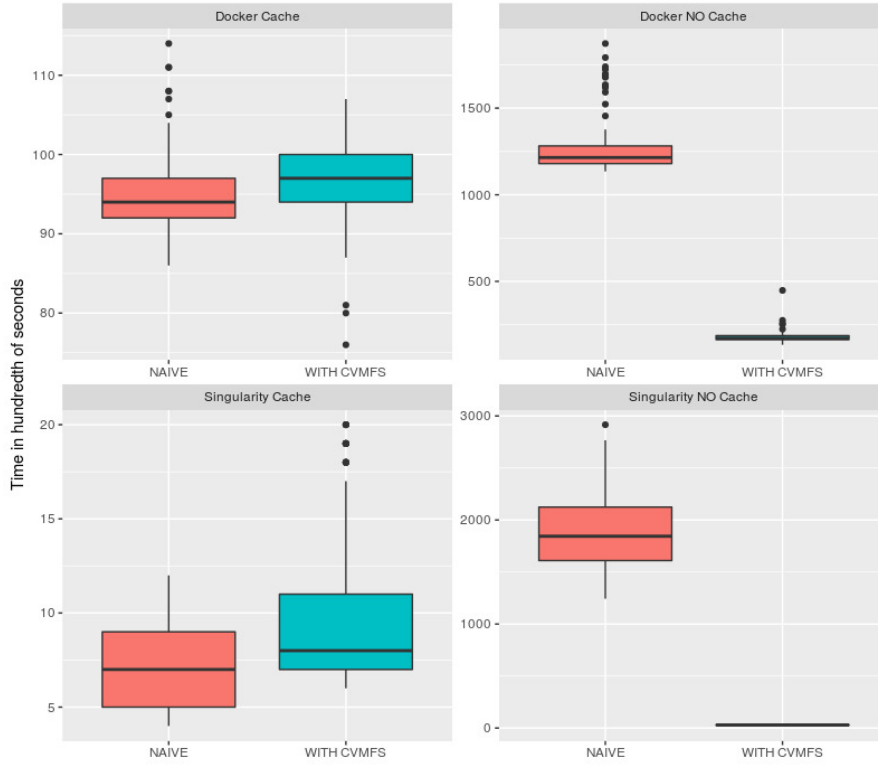


Figure 6.1: Startup times of the containers served using the proposed architecture and served using the standard Docker Registries.

	.flat	.layers
Size	28	27

Table 6.2: Apparent size in GB of the two folder *.layers* and *.flat*

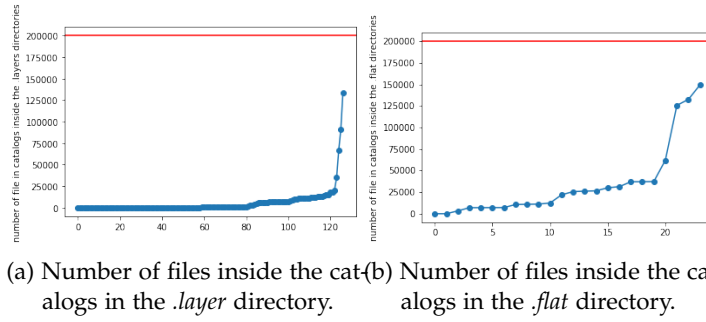
with the number of files in each catalog and the number of bytes that each catalog manage.

We can see that the Content-Addressable-Storage used by CVMFS helps a lot by reducing the amount of space required to host the repository by half, which make sense since the *.layers* and the *.flat* contains the exact same files.

FILE IN THE SUB-CATALOGS

CVMFS provides a rule of thumb about the number of files that should be hosted in each catalog to avoid stressing the sub-catalog too much. It suggest to limit the amount of file to less than 200000. Each request to the mounted repository needs to lookup values in at least one catalog, hence having small catalogs allow to keep each lookup fast enough.

Using again the CVMFS command `'cvmfs_server list-catalogs'` we obtain the number of files in each catalog.



We show the result in the graphs 6.2, 6.3, ?? and ??.

We can see how all the catalogs have less than 200000 entries respecting the suggestion of CVMFS. Moreover we can see that the number of entries in each catalog is not linearly distributed, few catalogs amount to a lot of entries while the majority of them have relatively few entries each.

COMPLEXITY

In order to estimate the complexity of managing the file-system we decide to use the cyclomatic complexity of the tool that create and manage the file-system itself that we discuss on Chapter 5. We used *gocyclo* [30] a command line tool that calculate the cyclomatic complexity of functions written in the Go(lang) language. We show the result of the cyclomatic analysis on Table 6.3.

While the cyclomatic complexity is very high is important to note that idiomatic golang codes requires to manually check every possible error returned by other functions, all these checks increase dramatically the cyclomatic complexity of the code. Indeed there are 19 error check without any logic in the code but simply returning early in the *ConvertWish* function.

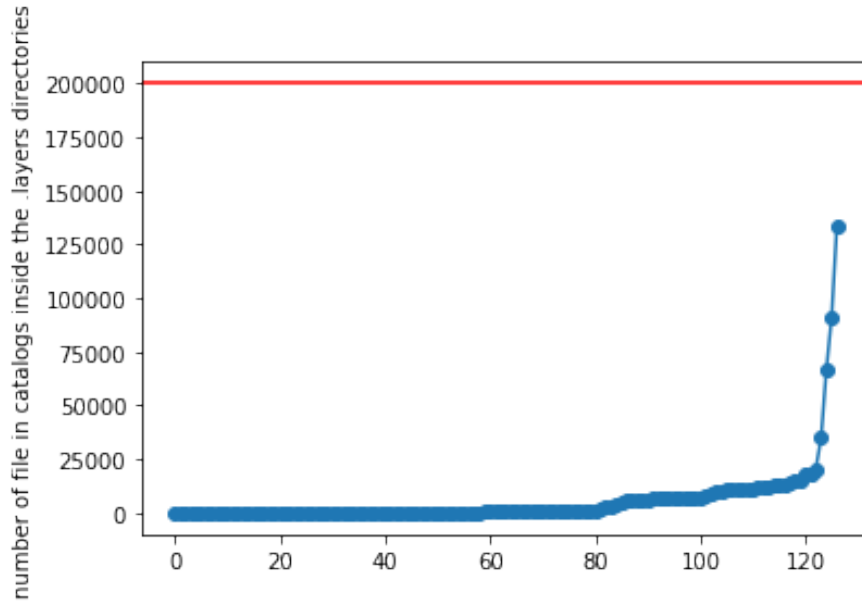


Figure 6.2: Number of files inside the catalogs in the *.layer* directory.

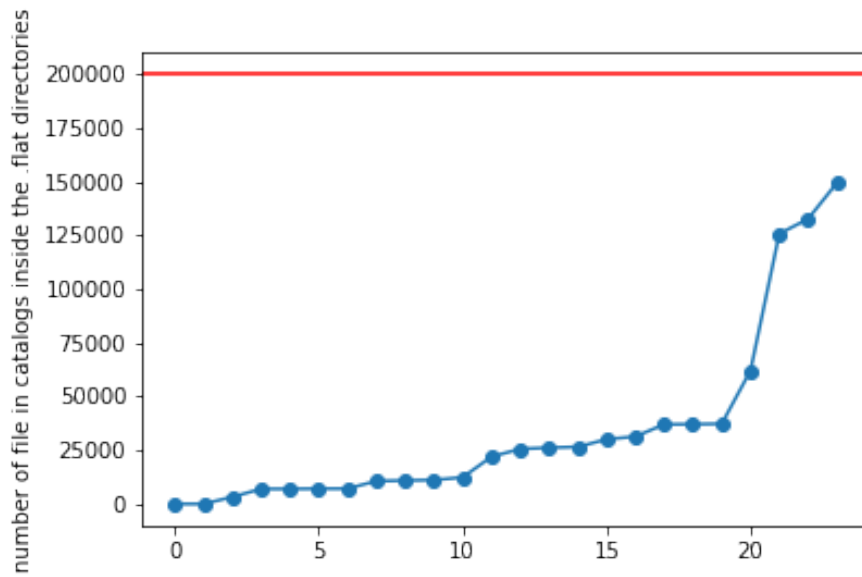


Figure 6.3: Number of files inside the catalogs in the *.flat* directory.

Cyclomatic Complexity	Function
37	ConvertWish
15	ParseImage
14	AddManifestToRemoveScheduler
12	GarbageCollectSingleLayer
12	(Image).GetChanges
12	(Image).downloadLayer
11	SaveLayersBacklink
11	IngestIntoCVMFS
9	requestAuthToken
9	CreateSymlinkIntoCVMFS
9	RemoveDirectory
8	(Image).GetLayers
7	FindImageToGarbageCollect
7	(Image).GetReference
6	AlreadyConverted
6	getBacklinkFromLayer
6	(*execCmd).Start

Table 6.3: Result of the cyclomatic complexity analysis, only function with complexity greater or equal to 3 are shown

FUTURE WORKS

1. New containers technologies
2. Improve ingestion of layers into the filesystem making everything in a single transaction
3. Decrease cyclomatic complexity

APPENDIX

CODE TO GET THE NUMBERS

Listing A.1: Script used to capture the startup time of singularity with image hostes in CVMFS using CVMFS cache

```
#!/bin/bash

for i in {1..101};
do
    /usr/bin/time -f "%U%S,%E" \
        singularity exec \
            /cvmfs/thin.osg.cern.ch/library/python:latest \
            python -c "quit()";
done
```

Listing A.2: Script used to capture the startup time of singularity with image hostes in CVMFS without cache

```
#!/bin/bash

for i in {1..101};
do
    #cleanup the cvmfs cache
    cvmfs_config wipecache >> /dev/null;

    /usr/bin/time -f "%U%S,%E" \
        singularity exec \
            /cvmfs/thin.osg.cern.ch/library/python:latest \
            python -c "quit()";
done
```

Listing A.3: Script used to capture the startup time of docker thin-images using both CVMFS and Docker cache

```
#!/bin/bash

docker pull thin-python:latest

for i in {1..101};
do
    /usr/bin/time -f "%U%S,%E" \
        docker run thin-python:latest python -c "quit()";
done
```

Listing A.4: Script used to capture the startup time of docker thin-images without Docker nor CVMFS cache

```
#!/bin/bash

for i in {1..101};
do
    #cleanup the cvmfs cache
    cvmfs_config wipecache >> /dev/null;

    # remove the layers from the docker cache
    docker rmi thin-python:latest -f >> /dev/null;

    # cleanup any remaining data from the system
    docker system prune -a -f;

    /usr/bin/time -f "%U,%S,%E" \
        docker run thin-python:latest python -c "quit()";
done
```

Listing A.5: Script used to capture the startup time of Docker standard images without Docker cache

```
#!/bin/bash

for i in {1..101};
do
    # remove the local image of python
    docker rmi -f python:latest;

    # cleanup any remaining data from the system
    docker system prune -a -f;

    /usr/bin/time -f "%S,%U,%E" \
        -o ~/start-time-docker-standard-no-cache.csv -a \
        docker run python:latest python -c "quit()";
done
```

Listing A.6: Script used to capture the startup time of Singularity running Docker standard images without Singularity cache

```
#!/bin/bash

for i in {1..101};
do
    # cleanup the singularity cache directory
    rm -rf /root/.singularity/docker;

    /usr/bin/time -f "%S,%U,%E" \
        -o ~/start-time-singularity-standard-no-cache.csv -a \
        singularity exec docker://python:latest \
        python -c "quit()";
done
```

Listing A.7: Script used to capture the startup time of Singularity running images unpacked on the local file-system, hence with cache

```
#!/bin/bash

for i in {1..101};
do
    /usr/bin/time -f "%S,%U,%E" \
        -o ~/start-time-singularity-standard.csv -a \
        singularity exec python_latest/ python -c "quit()";
done
```

Listing A.8: Script used to capture the startup time of docker thin-images without cache

```
#!/bin/bash
for i in {1..101};
do
    # clean the CVMFS cache
    cvmfs_config wipecache >> /dev/null;

    # clean the docker cache
    docker rmi thin-osg/library/python:latest -f >> /dev/null;

    # run and measure the start up time
    /usr/bin/time -f "%S,%U,%E" \
        -o ~/start-time-docker-no-cache.csv -a \
        docker run thin-osg/library/python:latest \
        python -c "quit()";
done
```


BIBLIOGRAPHY

- [1] Ali Anwar et al. “Improving Docker Registry Design Based on Production Workload Analysis.” In: *16th USENIX Conference on File and Storage Technologies (FAST 18)*. Oakland, CA: USENIX Association, 2018, pp. 265–278. ISBN: 978-1-931971-42-3. URL: <https://www.usenix.org/conference/fast18/presentation/anwar>.
- [2] J. O. Benson, J. J. Prevost, and P. Rad. “Survey of automated software deployment for computational and engineering research.” In: *2016 Annual IEEE Systems Conference (SysCon)*. 2016, pp. 1–6. DOI: [10.1109/SYSCON.2016.7490666](https://doi.org/10.1109/SYSCON.2016.7490666).
- [3] I Bird et al. *Update of the Computing Models of the WLCG and the LHC Experiments*. Tech. rep. CERN-LHCC-2014-014. LCG-TDR-002. 2014. URL: <http://cds.cern.ch/record/1695401>.
- [4] J. Blomer, P. Buncic, R. Meusel, G. Ganis, I. Sfiligoi, and D. Thain. “The Evolution of Global Scale Filesystems for Scientific Software Distribution.” In: *Computing in Science Engineering* 17.6 (2015), pp. 61–71. ISSN: 1521-9615. DOI: [10.1109/MCSE.2015.111](https://doi.org/10.1109/MCSE.2015.111).
- [5] Jakob Blomer, Predrag Buncic, and Thomas Fuhrmann. “Decentralized Data Storage and Processing in the Context of the LHC Experiments at CERN.” Presented 05 Jul 2012. 2011. URL: <https://cds.cern.ch/record/1462821>.
- [6] CERN. WLCG. URL: <https://web.archive.org/web/20181122162700/http://wlcg.web.cern.ch/> (visited on 11/22/2018).
- [7] G Compstellla, S Pagan Griso, D Lucchesi, I Sfiligoi, and D Thain. “CDF software distribution on the Grid using Parrot.” In: *Journal of Physics: Conference Series* 219.6 (2010), p. 062009. URL: <http://stacks.iop.org/1742-6596/219/i=6/a=062009>.
- [8] Christoph Eck et al. *LHC computing Grid: Technical Design Report. Version 1.06 (20 Jun 2005)*. Technical Design Report LCG. Geneva: CERN, 2005. URL: <http://cds.cern.ch/record/840543>.
- [9] The Linux Foundation. *Image Layer Filesystem Changeset*. URL: <https://web.archive.org/web/20181122170713/https://github.com/opencontainers/image-spec/blob/master/layer.md> (visited on 11/22/2018).
- [10] The Linux Foundation. *OCI Content Descriptors*. URL: <https://web.archive.org/web/20181122171756/https://github.com/opencontainers/image-spec/blob/master/descriptor.md> (visited on 11/22/2018).

- [11] The Linux Foundation. *Open Containers Initiative*. URL: <https://web.archive.org/web/20181122170509/https://www.opencontainers.org/> (visited on 11/22/2018).
- [12] N Hardi, J Blomer, G Ganis, and R Popescu. "Making containers lazy with Docker and CernVM-FS." In: *Journal of Physics: Conference Series* 1085.3 (2018), p. 032019. URL: <http://stacks.iop.org/1742-6596/1085/i=3/a=032019>.
- [13] Tyler Harter, Brandon Salmon, Rose Liu, Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau. "Slacker: Fast Distribution with Lazy Docker Containers." In: *14th USENIX Conference on File and Storage Technologies (FAST 16)*. Santa Clara, CA: USENIX Association, 2016, pp. 181–195. ISBN: 978-1-931971-28-7. URL: <https://www.usenix.org/conference/fast16/technical-sessions/presentation/harter>.
- [14] Docker Inc. *About storage drivers*. URL: <https://web.archive.org/web/20181122172916/https://docs.docker.com/storage/storagedriver/> (visited on 11/22/2018).
- [15] Docker Inc. *Docker Engine managed plugin system*. URL: <https://web.archive.org/web/20181122174052/https://docs.docker.com/engine/extend/> (visited on 11/22/2018).
- [16] Docker Inc. *Docker Registry*. URL: <https://web.archive.org/web/20181122172148/https://docs.docker.com/registry/> (visited on 11/22/2018).
- [17] Docker Inc. *Docker overview*. URL: <https://web.archive.org/web/20181122172634/https://docs.docker.com/engine/docker-overview/> (visited on 11/22/2018).
- [18] Docker Inc. *Docker tag*. URL: <https://web.archive.org/web/20181122172506/https://docs.docker.com/engine/reference/commandline/tag/> (visited on 11/22/2018).
- [19] Docker Inc. *What is a Container*. URL: <https://web.archive.org/web/20181122172321/https://www.docker.com/resources/what-container> (visited on 11/22/2018).
- [20] Red Hat. Inc. *Ansible is Simple It Automation*. URL: <https://web.archive.org/web/20181127163843/https://www.ansible.com/> (visited on 11/27/2018).
- [21] W. Kangjin, Y. Yong, L. Ying, L. Hanmei, and M. Lin. "FID: A Faster Image Distribution System for Docker Platform." In: *2017 IEEE 2nd International Workshops on Foundations and Applications of Self* Systems (FAS*W)*. 2017, pp. 191–198. DOI: [10.1109/FAS-W.2017.147](https://doi.org/10.1109/FAS-W.2017.147).
- [22] James H. Morris, Mahadev Satyanarayanan, Michael H. Conner, John H. Howard, David S. H. Rosenthal, and F. Donelson Smith. "Andrew: A Distributed Personal Computing Environment." In: *Commun. ACM* 29 (1986), pp. 184–201.

- [23] Simone Mosciatti. *repository-manager*, utility to manage the unpacked.cern.ch CVMFS repository. URL: <https://web.archive.org/web/20181125115146/https://github.com/cvmfs/docker-graphdriver/tree/devel/repository-manager> (visited on 11/25/2018).
- [24] Bogdan Nicolae, John Bresnahan, Kate Keahey, and Gabriel Antoniu. "Going Back and Forth: Efficient Multi-Deployment and Multi-Snapshotting on Clouds." In: *The 20th International ACM Symposium on High-Performance Parallel and Distributed Computing (HPDC 2011)*. San José, CA, United States, June 2011. URL: <https://hal.inria.fr/inria-00570682>.
- [25] Puppet. *Deliver better software, faster* | Puppet. URL: <https://web.archive.org/web/20181127164038/https://puppet.com/> (visited on 11/27/2018).
- [26] Nikolaus Rath. *libfuse*, The reference implementation of the Linux FUSE (Filesystem in Userspace) interface. URL: <https://web.archive.org/web/20181124093000/https://github.com/libfuse/libfuse> (visited on 11/24/2018).
- [27] Sylabs.io. *Singularity Appendix* | run. URL: <https://web.archive.org/web/20181122175328/https://www.sylabs.io/guides/2.6/user-guide/appendix.html#run> (visited on 11/22/2018).
- [28] Sylabs.io. *Singularity and Docker*. URL: https://web.archive.org/web/20181122174545/https://www.sylabs.io/guides/2.6/user-guide/singularity_and_docker.html (visited on 11/22/2018).
- [29] Sylabs.io. *Singularity*. URL: <https://web.archive.org/web/20181122154728/https://www.sylabs.io/> (visited on 11/22/2018).
- [30] fzippp. *gocycle*, Calculate cyclomatic complexities of functions in Go source code. URL: <https://web.archive.org/web/20181124194738/https://github.com/fzippp/gocyclo> (visited on 11/24/2018).

DECLARATION

Put your declaration here.

Saarbrücken, June 2018

André Miede & Ivo Pletikosić

COLOPHON

This document was typeset using the typographical look-and-feel `classicthesis` developed by André Miede and Ivo Pletikosić. The style was inspired by Robert Bringhurst's seminal book on typography "*The Elements of Typographic Style*". `classicthesis` is available for both L^AT_EX and L^YX:

<https://bitbucket.org/amiede/classicthesis/>

Happy users of `classicthesis` usually send a real postcard to the author, a collection of postcards received so far is featured here:

<http://postcards.miede.de/>

Thank you very much for your feedback and contribution.