

UNIVERSIDADE DA CORUÑA
SCHOOL OF INFORMATICS

Department of Electronics and Systems

Assessment, Design and Implementation of a Private Cloud for MapReduce Applications

Author: Marcos Salgueiro Balsa

Patricia González Gómez

Tutoring: Tomás Fernández Pena

José Carlos Cabaleiro Domínguez

Date: A Coruña, June 2013

*Give a man a fish, and you'll feed him for a day.
Teach a man to fish, and you'll feed him for a lifetime.*

Anne Isabella Thackeray Ritchie

*Great spirits have always encountered violent opposition from mediocre
minds.*

Albert Einstein

The supreme art of war is to subdue the enemy without fighting.

Sun Tzu, *The Art of War*

*[...] It takes these very simple-minded instructions – “Go fetch a number,
add it to this number, put the result there, perceive if it's greater than this
other number” – but executes them at a rate of, let's say, 1,000,000 per
second. At 1,000,000 per second, the results appear to be magic.*

Steven Paul Jobs

Summary

The history of computation has seen how the technology's unending evolution has promoted changes in its ways and means. Today, *tablets* and *smart-phones*, quantitatively inferiors managing and memorizing numbers, camp freely in a global market saturated with options. The tendency is clear: users will get to use more than one device to access the Internet and will like to have all of their data synchronized and at hand, all the time.

But that is only a part in the equation. At the other side of every service request there lays a server that must deal with an ever increasingly troubling traffic volume, while it maintains response delivery at outstanding delay times — low latency "may" have helped the infant Google rise above the competition. If we also added that the idea of surrounding every implementation effort with energetic efficiency is a transcendental requisite and not simply a good practice, we would have a perfect environment for the proliferation of new distributed paradigms as the *Cloud*. The Cloud is not an intrinsically new idea but an old concept abstraction: *virtualization*. The clouds' cornerstone is flexibility.

Another technology that is constantly making it to the headlines is *MapReduce*. If the Cloud centers around easing infrastructure exploitation, MapReduce's core strength lies in its speeding up driving large masses of unstructured data; with makes them an extraordinary computational tandem. This project puts forth a solution that allows for drawing on computational resources available exploiting both technologies together. Special emphasis has been placed in flexibility of access, being a web browser the only application

required to use the service; in simplifying the virtual cluster configuration, by including a self-managed minimum deployment; and in transparency and extensibility, by freeing source code and documentation as *OSS*, favoring its usage as starting point for larger installations.

Keywords

Distributed Computing, Virtualization, Cloud Computing, MapReduce, Open-Stack, Hadoop.

Contents

1	Abstract	1
1.1	Goals	3
1.2	Arrangement of the Document	4
2	Background	5
2.1	Cloud Computing	5
2.1.1	Architecture	7
2.1.2	Virtualization Techniques	10
2.1.3	Cloud IaaS frameworks	12
2.2	MapReduce Paradigm	13
2.2.1	Programming Model	13
2.2.2	Applicability of the Model	16
2.2.3	Processing Model	17
2.2.4	Fault Tolerance	19
2.2.5	Additional Characteristics	20
2.2.6	Other MapReduce Implementations	22
3	Experimental Assessment of IaaS Clouds	25
3.1	Assessment Methodology	25
3.2	Evaluated Frameworks	26
3.2.1	CloudStack	27
3.2.2	OpenNebula	29
3.2.3	OpenStack	30
3.3	Veredict	32

3.3.1	OpenStack Folsom	34
4	OpenStack Folsom	35
4.1	Global Architecture	35
4.2	Horizon	36
4.3	Keystone	39
4.4	Quantum	39
4.5	Compute	42
4.6	Glance	43
4.7	Storage	43
4.7.1	Cinder	44
4.7.2	Swift	44
5	Hadoop	47
5.1	The Beginnings	47
5.2	General Hadoop Architecture	48
5.3	Hadoop Distributed File System	50
5.3.1	Node Roles	51
5.3.2	Network Topology	53
5.4	Hadoop MapReduce	56
5.4.1	Node Roles	57
6	A Private Cloud for MapReduce Applications	61
6.1	Architecture	61
6.1.1	Design Diagrams	64
6.2	Implementación	66
6.2.1	Máquina virtual	67
6.2.2	Diagramas de implementación	69
	Bibliography	84

List of Figures

1.1	Demand in exponential growth. Source: <i>Cloudera Inc.</i>	2
1.2	Energy savings. Source: [1]	3
2.1	Layers in a cloud in production	6
2.2	Cloud Controller and Cloud Node	7
2.3	Cloud Controller in detail	9
2.4	Map function example (functional version)	14
2.5	Map function signature (MapReduce version)	14
2.6	Fold function example	15
2.7	Reduce function signature	15
2.8	MapReduce wordcount pseudocode. Source: [2]	16
2.9	MapReduce execution diagram. Source: [2]	18
3.1	General testing space	26
3.2	CloudStack 3.0.2 with XenServer hypervisor	28
3.3	OpenNebula 3.8 with KVM	30
3.4	Virtual OpenStack deployment	31
4.1	OpenStack Architecture	36
4.2	Example of an OpenStack Folsom deployment	37
4.3	Sequence Diagram — create instance	40
4.4	Virtual network deployment with Quantum	41
5.1	Hadoop over HDFS	49
5.2	Hadoop over another DFS	49

5.3	Typical HDFS deployment	51
5.4	Replica distance	55
5.5	Replication factor 3	57
5.6	Hadoop MapReduce with HDFS	57
6.1	Global Architecture	62
6.2	Detailed global architecture	63
6.3	Core qosh modular decomposition	63
6.4	Django setup	64
6.5	Use Cases Diagram	65
6.6	Web interface transitions	66
6.7	Diagrama de Clases del cliente de acceso REST (I)	67
6.8	Diagrama de Clases del cliente de acceso REST (II)	69
6.9	Diagrama de Clases — Django y Fabric	70
6.10	Diagrama de Clases — Objetos de Django	72
6.11	Diagrama de Clases — Formularios de Django	74
6.12	Diagrama Entidad-Relación	75
6.13	Diagrama de Secuencia (I)	76
6.14	Diagrama de Secuencia (II)	77

Chapter 1

Abstract

Over the last years there has been a continuous increase in the quantity of information generated with the Internet as the main driver. Furthermore, this information has reshaped from structured — and thus, susceptible to being expressed following a relational schema — to heterogeneous, which has kick-started the necessity to alter the way it is stored and transformed. As the figure 1.1 shows, those that were the undisputed back-end queens — relational database systems mostly — are seeing how their role is fading away due to their incapability to efficiently save unrelated heterogeneity.

As another related dimension, in the year 2000 many .com companies started upgrading their data centers to accommodate the inexorable demand peak that was going to follow. But it never came; and the bubble burst. What happened then was general underutilization — only 10% of Amazon's global computational resources were in use — that pushed the search for alternative means to export the surplus as a product. Amazon's own initiative unfolded in 2006 with the *AWS* (*Amazon Web Services*) appearance. AWS, among others, implements a public API for flexible on-demand infrastructure provisioning.

Since then, similar projects have proliferated generalizing how private clusters' unused computational capacity is to be serviced, trying to stay API-compatible with the AWS to facilitate interoperability and thus avoid client's

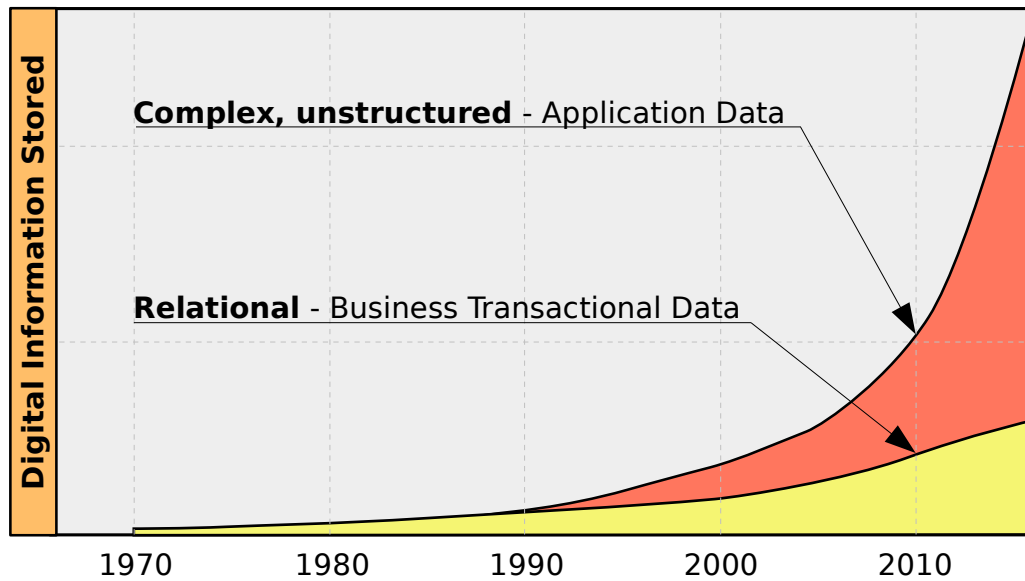


Figure 1.1: Demand in exponential growth. Source: Cloudera Inc.

swapping to more flexible providers.

Meanwhile, Google was also in the search for new mechanisms to exploit, with high performance and securely, their own private infrastructure to evolve the capability of their services. MapReduce, as a way to massively execute thousands transformations on input data, became a reality to thrust the generation of Google's humongous inverted index of the Internet [2]. Forthcoming contributions from Nutch's developers — by that time an Internet search engine prototype — to the MapReduce paradigm at *Yahoo!*, would traduce into the appearance of today's *de facto* standard in the field: Hadoop. Nowadays Hadoop is used in a myriad of backgrounds, ranging from travel booking sites to storing and servicing mobile data, ecommerce, image processing applications or searching for new forms of energy.

So, by stacking a MapReduce implementation atop elastic infrastructure an optimal exploitation of computational resources would be attainable, rapidly expanding or shrinking them on-demand, while simultaneously reducing the overall energy consumption required to accomplish processings

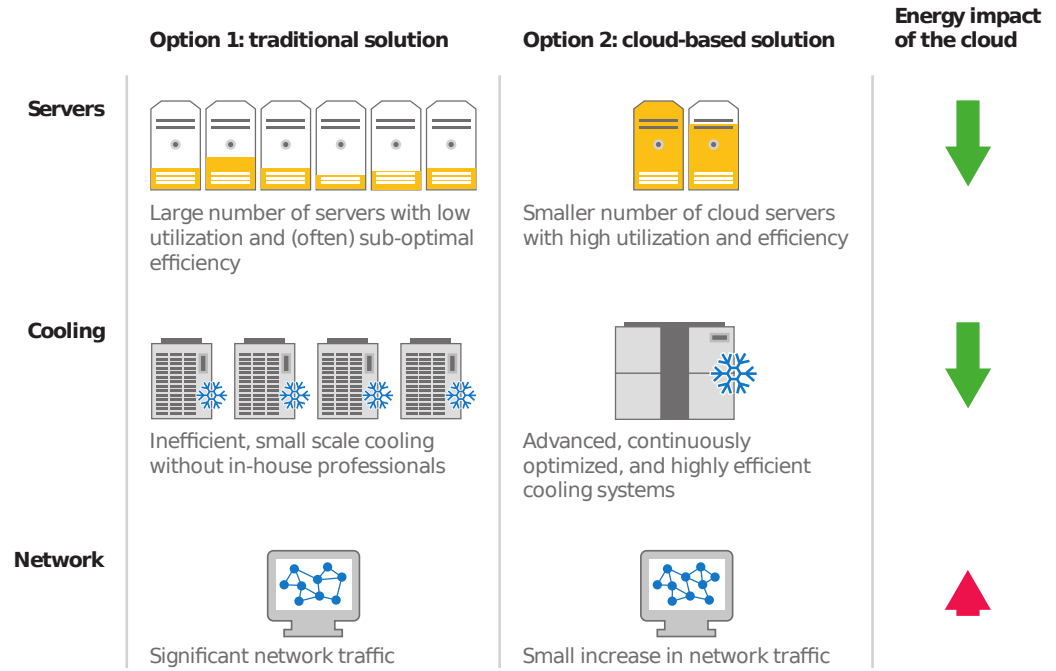


Figure 1.2: Energy savings. Source: [1]

(Figure 1.2).

1.1 Goals

The main goal with this project is to study the feasibility to develop a solution for a Cloud to drive MapReduce applications, with no need to know the particular Cloud structure and/or Hadoop configuration parameters.

In order to achieve such a simple execution model without compromising performance or applicability, a thorough analysis on different *IaaS Frameworks* will be carried out. Their features will be evaluated inside a virtual testing environment to finally narrow the selection to only one. Once an IaaS Framework had been chosen, the attention will be put towards choosing a MapReduce implementation to install over our virtual infrastructure.

Nonetheless, a mechanism to forward MapReduce execution requests will

be devised and implemented trying to focus on simplicity and universal access to this human-cloud-mapreduce interface. Yet, this transparency mustn't become an obstacle in exploiting the application or in fetching processed results. Privacy and security in communications and storage will be conveniently defined; we shan't forget it will be developed as a scaled-down model which could be infinitely scaled out.

1.2 Arrangement of the Document

The contents within this document are distributed as stated next. This first chapter introduces development guidelines in the abstract. Chapter 2 puts the reader closer to the fundamental Cloud Computing concepts — like its general architecture or virtualization —, along with the ones from the MapReduce paradigm. Chapter 3 describes an empirical evaluation of four private IaaS Cloud frameworks. Chapter 4 explores OpenStack Folsom's modular structure and particular inner workings. Analogically, chapter 5 unveils Hadoop's peculiarities as a MapReduce implementation.

The subsequent chapters center on detailing the project from diverse vantage points. Chapter 6 contains a series of design decisions and their accompanying UML diagrams. Chapter ?? gathers an analysis on performance in a real testing cluster. Chapter cap:conclusiones analyzes related papers highlighting how they compare to this solution. Finally, the main contributions of this project are discussed in addition to proposing future improvements to the implementation.

Two annexes have also been included. Annex ?? guides the reader throughout a quick single node installation. Annex ?? covers the definition of some of the concepts and technologies referred to in this text.

Chapter 2

Background

This second chapter tries to acquaint the reader with the key concepts that define Cloud Computing as well as the MapReduce archetype. Later successive elaborations to the project will lay on top of them.

2.1 Cloud Computing

In essence, Cloud Computing, or Cloud for short, is a distributed computing model that attempts to ease the consumption on-demand of that distributed infrastructure, by exporting it as virtual computational resources, platforms or services. However it may seem, the Cloud is no new technology but it introduces a new manner to exploit idle computing capacity. What it intends is to make orchestration of enormous data centers more flexible, so as to allow a user to start or destroy virtual machines as required — Infrastructure as a Service (*IaaS*) —, leverage a testing environment over a particular Operating System or software platform — Platform as a Service (*PaaS*) — or use a specific service like remote backup — Software as a Service (*SaaS*). Figure 2.1 shows the corresponding high level layer diagram of a generic Cloud.

Different IaaS frameworks will cover the functionality that is required to drive the cloud-defining *physical* infrastructure. Nonetheless, an effort to analyze, design, configure, install and maintain the intended service will

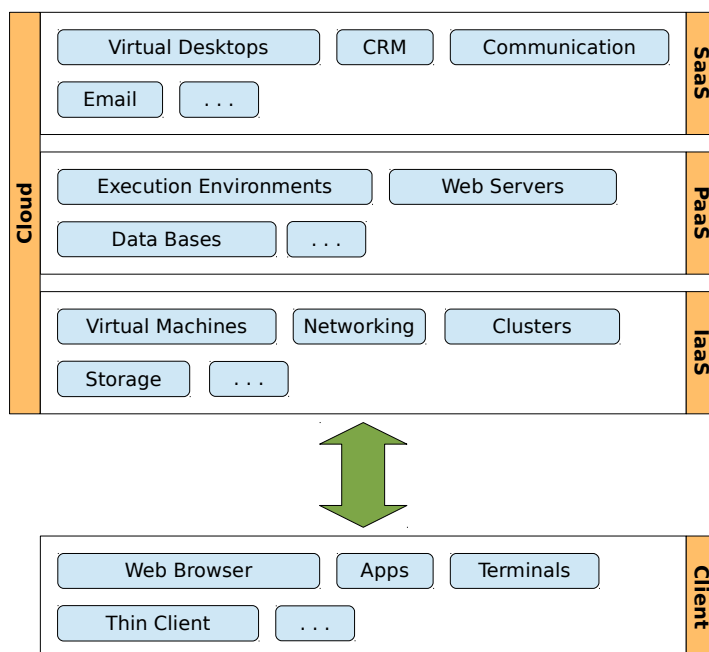


Figure 2.1: Layers in a cloud in production

be needed, bearing in mind that the degree of elaboration grows from IaaS services to SaaS ones. In effect, PaaS and SaaS layers are lied supported by those immediately under — software is implemented over a particular platform which, in turn, is also build upon a physical layer. Every Cloud Framework focuses on giving the option to configure a stable environment in which to run virtual machines defined by four variables: Virtual CPU count, virtual RAM, virtual persistent memory and virtual networking devices. Such an environment makes it possible to deploy virtual clusters upon which to install platforms or services to be subsequently consumed by users, bringing up the software layers that give form PaaS and SaaS paradigms respectively.

No less important cuestions like access control, execution permissions, quota or persistent or safe storage will also be present in all of the frameworks.

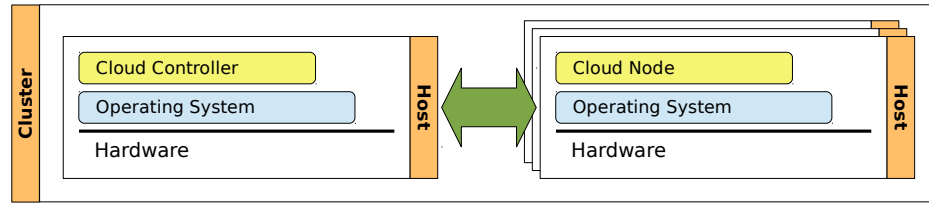


Figure 2.2: Cloud Controller and Cloud Node

2.1.1 Architecture

Figure 2.1 showed possible layers that could be found in a cloud deployment. Depending on the layers that are implemented, the particular framework and the role played by the cluster node, different particular modules will appear to make possible the consumption of configured services. These modules may be thought of as Cloud subsystems that connect each one of the parts that are required to execute virtual machines. Those virtual machines' capabilities are defined by the four variables previously discussed — VCPUS, RAM, HDD and networking. As there is no methodology dictating how those subsystems should be in terms of size and responsibility, and thus, each framework makes its own modular partition regarding infrastructure management.

Setting modularity apart, one common feature among different clouds is the separation of responsibility in two main roles: *Cloud Controller* and *Cloud Node*. Figure 2.2 shows a generic Cloud deployment in a cluster with both roles defined. The guidelines followed for having this two roles lies close to *Master-Slave* architectures' approach. In those, in the abstract, there's a set of computers labeled as coordinators which are expected to control execution, and another set made up with those machines that are to carry out the actual processing.

Within this general role distribution in a cluster, host computers or cluster nodes — labeled as Cloud Controllers or Cloud Nodes — cooperate in a synchronized fashion through *NTP* (*Network Time Protocol*) and communicate via message passing supported by asynchronous queues. To store

services' metadata and status they typically draw upon a *DBMS* (*Data Base Management System*) implementation, which is regularly kept running in a dedicated cluster node set sharded (distributed) between the members of the set.

Although there is no practical restriction to configuring both Cloud Controller and Cloud Node within a single computer in a cluster, this approach should be limited to development environments due to the considerable impact in performance that it would carry.

Cloud Controller

The fundamental task for a Controller is to maintain all of the cloud's constituent modules working together by coordinating their cooperation. As an example, it is a Controller's duty to:

- Authentication and authorization control.
- Available infrastructure resources recount.
- Quota management.
- Usage balance.
- User and project inventory.
- API exposure for service consumption.
- Real time cloud monitoring.

Being an essential part of a cloud as it is, the Controller node (not to be mistaken for the Cloud Node) is usually replicated in physically distinct computers. Figure 2.3 shows a Cloud Controller's architecture from a high level perspective.

As a general rule, clients will interact with clouds through a Web Service API — mostly *RESTful* APIs (*REpresentational State Transfer*). Those APIs vary slightly from company to vendor as usual, which forces clients to

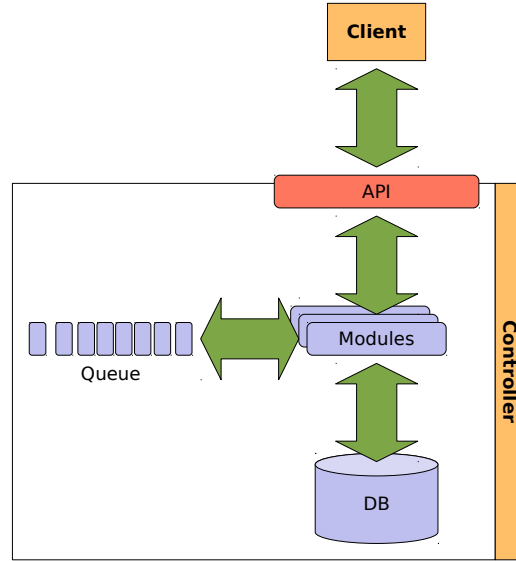


Figure 2.3: Cloud Controller in detail

be partially coupled to clouds. That is why there has been an increasing trend for unifying and standardizing those APIs in order to guarantee compatibility inter-framework. Of special mention is the cloud standard proposed by the *Open Grid Forum: OCCI (Open Cloud Computing Interface* [6]).

Cloud conforming modules support its functional requirements. Each one of them will have a well-defined responsibility, and so appear networking modules, access and security control modules, storage modules, etc. Many of them existed before the advent of Cloud Computing but they worked only locally. Inter-module communication is handled by means of an asynchronous message queue that guarantees an equally efficient broadcasting system outside of the Cloud Controller, i.e. the rest of the cluster nodes participating in the cloud.

To store and expose configuration data to the cluster in a single place while managing concurrent requests to update these data, every IaaS Framework evaluated resorts to a DBMS whose profiling must be properly tailored.

Hardware requirements on the cluster nodes vary from each particular framework implementation and the *QoS* expected, but, in the abstract, they normally need something around 10 GB of RAM, quad core CPU, Gigabit Ethernet and one TB of storage.

Cloud Node

If the Cloud Controller is entrusted the cloud's correct functioning acting like a glue for its parts, the actual task processing is performed in the Cloud Nodes; that is, the VCPU, VRAM, VHDD are going to be mapped from the corresponding CPU, RAM and HDD from the real nodes of the cluster.

Cloud Nodes may be heterogeneous according to their hardware characteristics. They will configure a resource set that, seen from the outside of the cluster, will appear to be a homogeneous whole where the summation of capacities of every participating node is the cloud's dimension. Further, this homogeneous space could be provisioned, as discussed above, on demand. It is the Cloud Controller's responsibility single out the optimal distribution of virtual servers throughout the cluster, attending to the physical aspects of both the virtual machine and the computer in which the former will run.

The most important subsystem in a Cloud Controller is the *hypervisor* or *VMM* (*Virtual Machine Monitor*). The hypervisor is responsible for making possible the execution of virtual servers — or virtual instances following the AWS nomenclature — by creating the virtual architecture needed and a *virtual execution domain* managed with the help of the operating system kernel. To generate this architecture there fundamentally exist three techniques: *Emulation*, *Paravirtualization* and *Hardware Virtualization* or *Full Virtualization*. Different hypervisors will support them in a different degree, but most will cover only one of them.

2.1.2 Virtualization Techniques

What follows is a brief review of the main methods to create virtual infrastructure.

Emulation

Emulation is the most general virtualization method, in a sense that it does not call for anything special be present in the underlying hardware. However, it also carries the highest penalization in terms of performance. With emulation, every structure sustaining the virtual machine operation is created as a functional software copy of its hardware counterpart; i. e., every machine instruction to be executed in the virtual hardware must be run software-wise first, and then be translated on the fly into another machine instruction runnable in the physical domain — the cluster node. The interpreter implementation and the divergence between emulated and real hardware will directly impact the translation overhead. This fact hinders the emulation from being widely employed in performance-critical deployments. Nonetheless, thanks to its operating flexibility it's generally used as a mechanism to support legacy systems. Besides, the kernel in the guest operating system — the kernel in the virtual machines's — operating system needs no alteration whatsoever, and the cluster node's kernel need only load a module.

Hardware Virtualization

Hardware Virtualization, on the contrary, allows host's processes to run directly atop the physical hardware layer, with no interpretation. Logically, this provides a considerable speedup from emulation, though imposes a special treatment to be given to its virtual processes. Regarding CPUs, both AMD's and Intel's support virtual process execution — which is the capacity to run processes belonging to the virtual domain with little performance overhead — as far as the convenient hardware extensions are present (*SVM* and *VT-x* respectively [7]). Just as what happened with emulation, an unaltered host's kernel may be used. This fact is of relative importance as if wasn't so it would limit the myriad of OSs that could be installed as guests. Lastly, it should be pointed out that the hardware architecture is exposed to the VM as it is, i. e. with no software middleware.

Paravirtualization

Paravirtualization uses a different approach. To begin with, it is indispensable that the guest's kernel be modified to make it capable of interacting with a paravirtualized environment. When the guest runs, the hypervisor will separate those regions of instructions that have to be executed in kernel mode in the CPU, from those in user mode which will be executed as regular host processes. Subsequently, the hypervisor will manage an on-contract execution between host and guest allowing the latter to run kernel mode sections as if pertaining to the real execution domain — as if they were processes running in the host, not in the guest — with almost no performance slowdown. Paravirtualization, in turn, does not require an special hardware extension be present.

2.1.3 Cloud IaaS frameworks

Cloud IaaS frameworks are those software systems managing the abstraction of complexity associated with on demand provisioning and administering failure-prone generic infrastructure. In spite of being almost all of them open sourced — which fosters reusability and collaboration —, they have evolved in different frames. This fact has raised a condition of lacking outwards interoperability, maturing non-standard APIs; though today those divergences are fading away. These frameworks and APIs are product of the efforts to improve and ease controlling the underlying particular clusters on which they germinated. Thus, it is no surprising their advances had originated parallely with the infrastructure they drove, leaving compatibility in the background.

Slowly but steadily these managing systems became larger in reach and responsibility boosted by an increasingly interest in the sector. In the end, it happened that software and systems engineering made them more abstract, so they finally overlapped functionally. AWS appearance finished forging the latent standardization need, and thus, as of today, most frameworks offer APIs closer and closer to Amazon's — nowadays the de-facto standard —

and OCCI's [6].

2.2 MapReduce Paradigm

The origin of the paradigm centers around a paper publication of two Google employees [2]. In this paper they explained a method implementation devised to abstract the common parts present in distributed computing that rendered simple but large problems much more complex to solve when paralleling their execution on massive clusters. A concise definition states that MapReduce is “*a data processing model and execution environment that runs on large clusters of commodity computers*” [8].

2.2.1 Programming Model

The MapReduce programming model requires the developer express his problem as a partition of two well-defined pieces. A first part deals with the reading of input data and with producing a set of intermediate results that will be scattered over the cluster nodes. These intermediate transformations will be grouped according to an intermediate key value. A second phase begins with that grouping of intermediate results and concludes when every *reduce* operation on the groupings succeeds. Seen from another vantage point, the first phase corresponds, broadly speaking, to the behavior of the functional *map* and the second to the functional *fold*.

In terms of the MapReduce model, these functional paradigm concepts give rise to *Map* and *Reduce* functions. Both Map and Reduce have to be supplied by the developer, which may force a deviation in breaking the original problem down. As counterpart, the MapReduce model will deal with parallelizing the computation, distributing input data across the cluster, handling exceptions that could raise and recovering output results; everything transparent to the programmer.

$$\begin{aligned} \mathbf{map} &: (\alpha \rightarrow \beta) \rightarrow \alpha \text{ list} \rightarrow \beta \text{ list} \\ \mathbf{map} \ (\mathbf{pow} \ 2) \ [1, 2, 3] &\Rightarrow [1, 4, 9] \end{aligned}$$

Figure 2.4: Map function example (functional version)

$$\begin{aligned} \mathbf{map} &: (k1, v1) \rightarrow (k2, v2) \text{ list} \\ \mathbf{k} &: \text{clave} \\ \mathbf{v} &: \text{valor} \\ (\mathbf{kn}, \mathbf{vn}) &: \text{par } (\text{clave}, \text{valor}) \text{ en un dominio } n \end{aligned}$$

Figure 2.5: Map function signature (MapReduce version)

Función Map

The typical functional map takes any function F and a list of elements L or, in general, any recursive data structure, to return a list resulting from applying F to each element of L . Figure 2.4 shows its signature and an example.

In its MapReduce realization, map function receives a tuple as input and produces another tuple $(key, value)$ as intermediate output. It is the MapReduce library who is responsible for feeding the map function by mutating the data contained in input files into $(key, value)$ pairs. Then, it deals with grouping those intermediate tuples by key before passing them in as input to the reduce function. Input and output data types correspond to those shown in the function signature figure 2.5.

Reduce function

The typical functional fold expects any function G , a list L , or generally any type of recursive data structure, and any initial element I , subtype of L 's elements. Fold returns the value in I resulting from building up the intermediate values generated after applying G to each element in L . Figure 2.6 presents fold signature as well as an example.

$$\mathbf{fold} : (\alpha \rightarrow \beta \rightarrow \alpha) \rightarrow \alpha \rightarrow \beta \text{ list} \rightarrow \alpha$$

$$\mathbf{fold} (+) 0 [1, 2, 3] \Rightarrow 6$$

Figure 2.6: Fold function example

$$\mathbf{reduce} : (k2, v2 \text{ list}) \rightarrow v2 \text{ list}$$

$$\mathbf{k} : \text{clave}$$

$$\mathbf{v} : \text{valor}$$

$$(\mathbf{kn}, \mathbf{vn}) : \text{par } (\text{clave}, \text{valor}) \text{ en un dominio } n$$

Figure 2.7: Reduce function signature

Contrary to map, reduce expects the intermediate groups as input to produce a smaller set of values for each group as output, because reduce will iteratively *fold* the groupings into values. Those reduced intermediate values will be passed in again to the reduce function if more values with the same key appeared from subsequent maps. Reduce signature is shown on figure 2.7. Just as happens with map, MapReduce handles the transmission of intermediate results out from map into reduce. The model also describes the possibility to define a *Combiner* function that would act after map partially reducing the values within the same grouping to lower network traffic — the combiner usually runs in the same machine as the map.

A word counter in MapReduce

As an example, figure 2.8 shows the pseudocode of a MapReduce application to count the number of words in a document set.

In a wordcount execution flow the following is going to happen: map is going to be presented with a set of names containing all of the documents in plain text whose words will be counted. Map will subsequently iterate over each document in the set emitting the tuple $(\langle \text{word} \rangle, "1")$ for each word found. Thus, an explosion of intermediate pairs will be generated as

```
Map (String key, String value):  
  // key:  documento name  
  // value:  document contents  
  for each word w in value:  
    EmitIntermediate (w, "1");  
  
Reduce (String key, Iterator values):  
  // key:  a word  
  // values:  an Iterable over intermediate counts of the word key  
  int result = 0;  
  for each v in values:  
    Emit (AsString (result));
```

Figure 2.8: MapReduce wordcount pseudocode. Source: [2]

output of map, will be distributed over the network and progressively folded in the reduce phase. Reduce is going to be input every pair generated by map but under a different form. Reduce will accept on each invocation the pair ($\langle word \rangle$, $list("1")$). The list of "1"s, or generically an `Iterable` over "1"s, will contain as many elements as instances of the word $\langle word \rangle$ there were in the document set — this supposing that the map phase were over before starting the reduce phase and that every word $\langle word \rangle$ were submitted to the same reducer in the cluster — a cluster node executing the reduce function.

Once the flow had been completed, MapReduce would return a listing with every word in the documents and the number of times it appeared.

2.2.2 Applicability of the Model

The myriad of problems that could be expressed following the MapReduce programming paradigm is clearly reflected in [2], a subset of them being:

- Distributed grep: Map emits every line matching the regular expres-

sion. Reduce only forwards its input to its output acting as identity function.

- Count of URL access frequency: Like wordcount.
- Reverse web-link graph: For each URL contained in a web document, map generates the pair $(\langle target_URL \rangle, \langle source_URL \rangle)$. Reduce will emit the pair $(target, list(source))$.
- Inverted index: Map parses each document and emits a series of tuples in the form $(\langle word \rangle, \langle document_id \rangle)$. All of them are passed as input to reduce that generates the sequence of pairs $(\langle word \rangle, list(document_id))$.

2.2.3 Processing Model

Besides defining the structure that the applications willing to leverage the MapReduce capabilities will have to follow — so that they need not code their own distribution mechanisms —, with [2] an implementation of the model was introduced which allowed Google to stay protocol, architecture and system agnostic while keeping their commodity clusters on full utilization. This agnosticism allows for deploying vendor-lock-free distributed systems.

The MapReduce model works by receiving self-contained processing requests called *jobs*. Each job is a *partition* of smaller duties called *tasks*. A job won't be completed until no task is pending for finishing execution. The processing model main intent is to distribute the tasks throughout the cluster in a way that reduced job latency. In general, it can be stated that task processing on each phase is done in parallel and phases execute in sequence; yet, it is not needed for reduce to wait until map is complete.

Figure 2.9 shows a summary of a typical execution flow. It is interesting enough to deepen in its details as many other MapReduce implementations will present similar approaches.

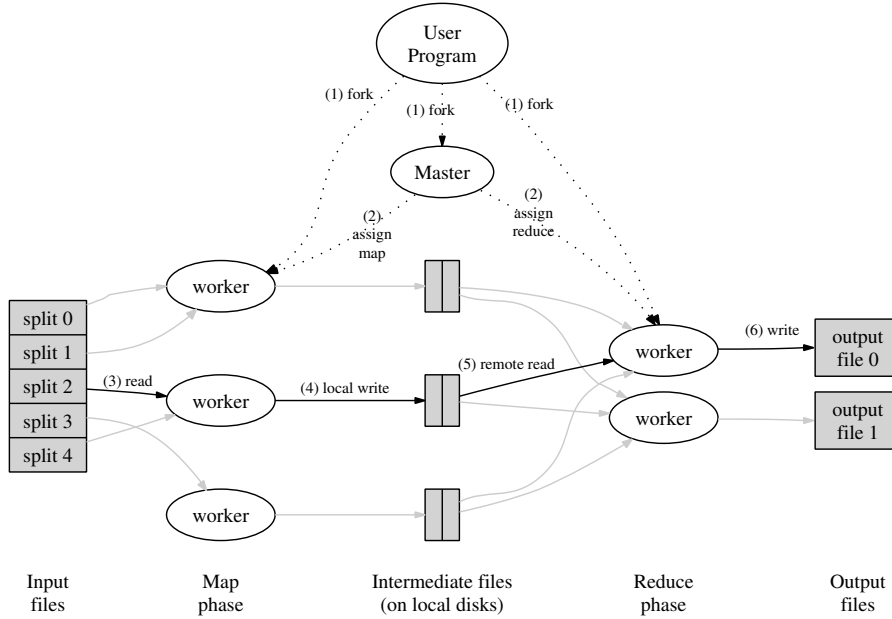


Figure 2.9: MapReduce execution diagram. Source: [2]

1. MapReduce divides input files in M parts, the size of which is parameterized, and distributes as many copies of the MapReduce user algorithm as nodes participate in the computation.
2. From this moment each program copy resides in a cluster node. A random copy is chosen among them and labeled as the *Master Replica*, effectively assigning the *Master Role* to the node holding the replica; every other node in the cluster is designated with the *Worker Role*. Those worker nodes will receive the actual MapReduce tasks and their execution will be driven from the master node. There will be M map tasks and R reduce tasks.
3. Workers assigned with map tasks read their corresponding portions of the input files and parse the contents generating tuples $(key, value)$ that will be submitted to the map function for processing. Map outputs are stored in memory as a cache.

4. Periodically, those pairs in memory are dumped to a local disk — dumped to a drive of the node that is executing the map function — and partitioned into R regions. Their path on disk is then sent back to the master, responsible for forwarding these paths to *reduce workers* or *reducers*.
5. Now, when a reducer is notified that it should start processing, the path to the data of the reduction is sent along and the reducer will fetch them directly from the mapper via *RPC* (*Remote Procedure Call*). Before actually invoking the reduce function, the node itself will sort the intermediate pairs by key.
6. Lastly, the reducer iterates over the key-sorted pairs submitting to the user-defined reduce function the key and the `Iterable` of values associated to the key. The output of the reduce function for the reducer partition is appended to a file stored over the distributed file system.

When every map and reduce tasks had succeeded, the partitioned output space — the file set within each partition — would be returned back to the client application that had made the MapReduce invocation.

This processing model is abstract enough as to be employed to the resolution of indeterminately large problems running on huge clusters.

2.2.4 Fault Tolerance

The idea of providing an environment to execute jobs long enough to require large sets of computing machines to keep the latency within reasonable timings, calls for the definition of a policy able to assure a degree of tolerance to failure. If unattended, those failures would lead to errors; some would cause finished tasks to get lost, others would put intermediate data offline. Consequently, if no measures were taken to prevent or deal with failure, job throughput would humble as some would have to be rescheduled all along.

The MapReduce model describes a policy foreseeing a series within an execution flow and duly implements a series of actions against them.

Worker Failure

The least taxing of the problems. To control that every worker is up, the master node pings them periodically. If a worker did not reply to pings repeatedly, it would be marked as failed.

A worker marked failed will neither be scheduled new tasks nor will be remotely accessed by reducers to load intermediate map results that it may had; a fact that could prevent the workflow from succeeding. If so were the case, the access to these data would be resolved by the master labeling the results of the failing tasks as *idle*, so that they could be rescheduled a later time to store the results in an active worker.

Master Failure

Failure of a master node is more troublesome. The proposed approach consists in making the master periodically create a snapshot from which to restore to a previous state if it went unexpectedly down. It is a harder problem than a worker failure mainly because there can only be one master per cluster, and the time it would take another node to take over the master role would leave the scheduling pipeline stalled. The master being in a single machine has also the benefit of lowering the probability of failure, precisely why in the original paper [2] it had been put forward that the entire job be canceled. Still, as there is no good design to leave a *single point of failure*, subsequent MapReduce implementations have proposed to replicate the master in other nodes in the same cluster.

2.2.5 Additional Characteristics

What follows is a summary of additional features of the original MapReduce implementation.

Locality

The typical bottleneck in a modern deployment is network bandwidth. In MapReduce executions, the information flows into the cluster from the external client. As already discussed, each node in a MapReduce cluster holds a certain amount of the input data and shares its processing capacity to be used for particular MapReduce tasks over those data. Each stage in the MapReduce executing pipeline requires a lot of traffic to be handled by the network which would reduce throughput if no wide enough channel were deployed nor a locality exploiting strategy were implemented.

In fact, MapReduce explores a method to use locality as an additional resource. The idea is for the distributed file system to place data as close as possible to where they will be transformed — it will try to store data in the mappers' and reducers' local drives —, effectively diminishing the transport over the net.

Complexity

A priori, variables M and R , the number of partitions of the input space and of the intermediate space respectively, may be configured to take any value whatsoever. Yet, there exist certain practical limits to their values. For every running job the master will have to make $O(M + R)$ scheduling decisions — if no error forced the master to reschedule tasks —, as each partition of the input space will have to be submitted to a mapper and each intermediate partition will have to be transmitted to a reducer, coming to $O(M + R)$ as the expression of *temporal complexity*. Regarding *spatial complexity*, the master will have to maintain $O(M \cdot R)$ as piece of state in memory as the intermediate results of a map task may be propagated to every piece R of the reduce space.

Backup Tasks

A situation could arise in which a cluster node be executing map or reduce tasks much slower than it theoretically could. Such a circumstance may arise with a damaged drive which would cause read and write operations to slow down. And since jobs complete when all of its composing tasks had been finished, the faulted node (*the straggler*) would be curbing the global throughput. To alleviate this handicap, when few tasks are left incomplete for a particular job, *Backup Tasks* are created and submitted to additional workers, making a single task be executed twice concurrently. By the time one copy of the task succeeds it will be labeled completed, duly reducing the impact of stragglers at the cost of wasting computational resources.

Combiner Function

Many times it happens that there exists a good number of repeated intermediate pairs. Taking wordcount as an example, it can be easily seen that every mapper will generate as many tuples (*"a", "1"*) as *a*'s there are in the input documents. A mechanism to lower the tuples that will have to be emitted to reducers is to allow for the definition of a *Combiner Function* to group outputs from the map function — and in the same mapper node — before sending them out over the network, effectively cutting down traffic.

In fact, it is usual for both combiner and reduce functions to share the same implementation, even though the former writes its output to local disk while the latter writes directly to the distributed file system.

2.2.6 Other MapReduce Implementations

Since 2004 multiple frameworks that implement the ideas exposed in the paper [2] have been coming out. The next listing clearly shows the impact MapReduce has created.

Hadoop [8] One of the first implementations to cover the MapReduce processing model and framework of reference to other MapReduce codifi-

cations. It is by far the most widely deployed, tested, configured and profiled today.

GridGain [9] Commercial and centered around in-memory processing to speedup execution: lower data access latency at the expense of smaller I/O space.

Twister [10] Developed as a research project of the University of Indiana, tries to separate and abstract the common parts required to run MapReduce workflows in order to keep them longer in the cluster's distributed memory. With such an approach, the time taken to configure mappers and reducers in multiple executions is lowered by doing their set up only once. Thus, *Twister* really shines in executing *iterative* MapReduce jobs — those jobs where maps and reduces do not happen in sequence once, but need instead a multitude of complete map-reduce cycles to succeed.

GATK [11] Used for genetic research to sequence and evaluate DNA fragments from multiple species.

Qizmt [12] Written in C# and deployed in MySpace.

misco [13] Written 100% Python and based on previous work at Nokia it is posed as a MapReduce implementation capable of running in mobile devices.

Peregrine [14] By optimizing how intermediate results are transformed and by passing every I/O operation throughout an asynchronous queue, its developers claim to have formidably accelerated task execution rate.

Mars [15] Implemented in *NVIDIA CUDA*, it revolves around extracting higher performance by moving the map and reduce operations into the graphic card. It is supposed to improve processing throughput by over an order of magnitude.

Hadoop is undoubtedly the most used MapReduce implementation nowadays. Its open source nature and its flexibility, both for processing and storing, have been reporting back an increasing interest from the IT industry. This has brought out many pluggable extensions that enhance Hadoop's applicability.

Chapter 3

Experimental Assessment of IaaS Clouds

In this chapter we will be reviewing the most used frameworks to drive IaaS Clouds. An initial selection will be made and it will be progressively shrunk following certain criteria like maturity, ease of use or documentation quality, until one remains. A deep study will be carried out on that prevailing one.

3.1 Assessment Methodology

A thorough evaluation of the capabilities of the different frameworks is not possible unless an actual deployment is carried through. The virtual infrastructure that is generated when a cloud has finished installing, no matter how small the deployment, is large and complex. Besides, trying to *emulate* the real hardware that will support the cloud is meager at times, e.g. if full hardware virtualization were used, the hypervisor would have to be allowed direct access to the CPU. *Nested Hardware Virtualization* — the capacity for a CPU to export its native virtualization capabilities to a guest running atop a host node, or the ability to use full hardware virtualization *inside* a virtual machine —, does not currently enjoy widespread adoption as it requires implementation efforts from both CPU designers and virtualization software

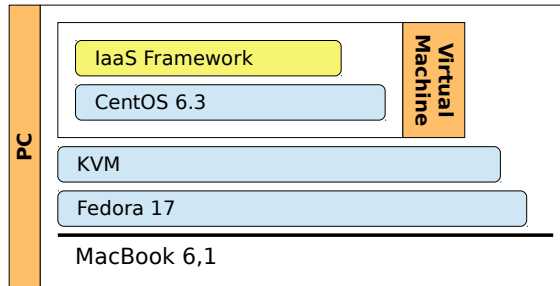


Figure 3.1: General testing space

developers. This means that it will not be possible for us to fully appraise the myriad IaaS Cloud solutions by creating a virtual cluster over which to deploy our clouds, and make performance measures.

To diagnose the superiority of one of them over the rest, a scaled down setup will be completed to evaluate the proficiency in maintaining the IaaS service running in spite of the reduced infrastructure. Quality and transparency in the documentation, as well as community support and engagement will also be born in mind.

The testing environment follows the organization shown in figure 3.1.

3.2 Evaluated Frameworks

The frameworks are:

- CloudStack
- Eucalyptus
- OpenNebula
- OpenStack

From an exclusively functional vantage point, the four of them cover clearly the requirements imposed with the project, which, in short, it would

be the faculty to run and manage the lifecycle of an indeterminate number of custom VMs, tailored for MapReduce executions, through an API that would allow for the definition of a simple job control interface.

Eucalyptus and *OpenStack* take a more modular approach to the solution, unveiling smaller functional parts, while at the same time decoupling those modules. With a module set in this fashion, installations become more flexible and tougher on par. However, to contain the operative effort, OpenStack ships with a series of scripts that help managing its deployments. Regarding their system requirements, they all support installations in modern Linux distributions with KVM or Xen as hypervisors. When dealing with a real deployment, the framework of choice will likely depend more on the existing platform than on particular limitations that any cloud may have.

As a side note, it is remarkable the lack of interoperability among them. All of them try to adhere to the AWS API in different degree — some of them partially support it, others use *adaptors* to it. OpenNebula, OpenStack and Eucalyptus have demonstrated to be carrying on coding efforts to fully support the OGF’s standard: OCCI.

Eucalyptus, in spite of being the first to fully cover the AWS API, which is merely anecdotal nowadays, has two obstacles that hinder its evaluation. First, it is not fully open sourced: **VMware Broker** which brings the opportunity to use virtualized infrastructure based on VMware technology, is only available to paid subscribers. And second, it is impossible to setup Eucalyptus within a VM to test start-up time or installation complexity, for example, as it explains its installation guide [16]. Both limitations make Eucalyptus back out from the evaluation list. The rest have been compared after their set up and configuration.

3.2.1 CloudStack

CloudStack installation guide ([17]) describes the series of steps that a systems engineer should follow in order to complete a minimum CloudStack deployment. It clearly determines that Cloud Nodes’ CPUs have to support

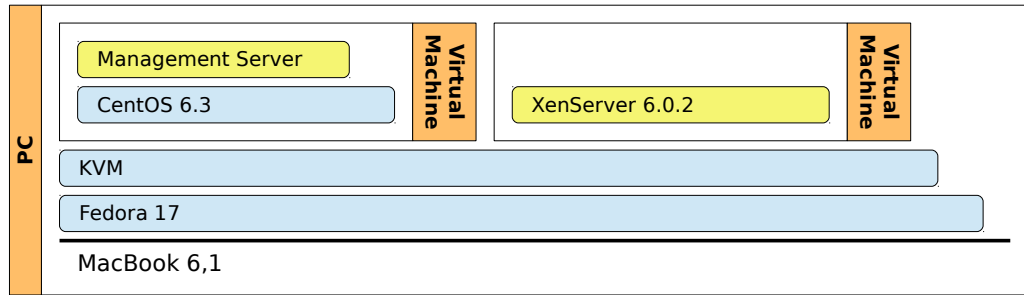


Figure 3.2: CloudStack 3.0.2 with XenServer hypervisor

virtualization extensions for CloudStack to start Xen or KVM-based VMs. Which happens to be a similar limitation to Eucalyptus'. However, the fact that CloudStack would become part of *The Apache Software Foundation* from version 4 onward ([18]), and the reality of a Citrix technical article opening the door to CloudStack deployments over Cloud Nodes lacking virtualization extensions ([20]), made us arrange the layout shown in figure 3.2.

Following the advanced and quick installation guides — [17] and [21] — the process may be summarized in the steps below:

- Two VMs were created to contain CloudStack *Management Server* (*MS*) and XenServer hypervisor: 1 GB of RAM for the MS, 3 GB of RAM for Xen, 20 GB HDD, *ACPI* and *APIC* for both.
- For the MS:
 - *CentOS 6.3* was downloaded, installed and *yum-updated*.
 - The VM was named *cloudstack*.
 - Likewise, a user named *cloudstack* was registered and added to the *sudoers* list.
 - The quick installation guide was followed to conclude the process.
- For Xen:
 - XenServer 6.0.2 was downloaded from Citrix web site.

- The notes contained in the quick installation guide and in the XenServer configuration manual [22] were followed to perform the configuration.
- Additionally:
 - Before specifying the execution environment, defining the cluster, primary and secondary storage, etc. a global flag had to be set to permit nodes with no virtualization extensions [19].
 - Once the configured infrastructure was online:
 - * The CentOS 6.3 image was uploaded to the MS.
 - * A `SimpleHTTPServer` — a Python micro HTTP server — was started on port `443` in the MS.
 - * A rule was added in `iptables` to let traffic through on port `443` in the MS.
 - * The image was loaded to the cloud from the web interface.

3.2.2 OpenNebula

If balanced against the installation procedure just described, the effort for setting up OpenNebula 3.8 is lighter. The process that has been followed to configure the OpenNebula deployment contained in figure 3.3 stems directly from [23], the official installation guide. The subsequent steps serve as summary to the process.

- A supporting VM was created to fully contain OpenNebula: 1 GB of RAM, 8 GB for the HDD, ACPI and APIC.
- CentOS 6.3 was downloaded, installed and yum-updated.
- The VM was named *opennebula*.
- A user named *opennebula* was also registered and added to the sudoers list.

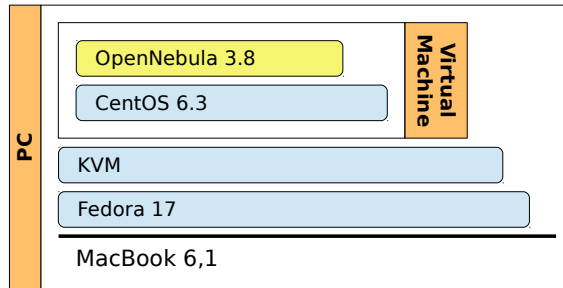


Figure 3.3: OpenNebula 3.8 with KVM

- SELinux was stopped.
- The configuration guide ([23]) was followed to complete the process. Afterwards, the next series of actions were carried out to test the framework:
 - By default, OpenNebula’s web interface module (`sunstone`) attaches to `lo`. In order to interact with `sunstone` from outside of the containing VM, it was necessary to change the configuration file (`/etc/one`) so that `sunstone` attached to the external networking interface `eth0`. The very same happened with `occi`, the REST service that exposes the cloud API.
 - Furthermore, `iptables` was modified for letting through traffic on port `9869`; where `sunstone` listens by default.

3.2.3 OpenStack

OpenStack case is striking for many reasons. It represents the convergence of two different needs: the computationally-driven in NASA and the storage-bound in Rackspace. Complementarily, both Red Hat and Canonical had shown their interest in the platform by collaborating with their scripts, deploying utilities and cloud-optimized images.

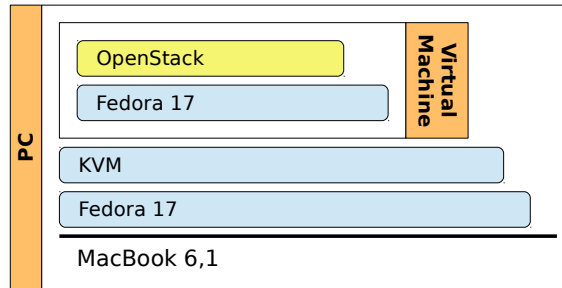


Figure 3.4: Virtual OpenStack deployment

Just as with OpenNebula, the complete execution environment was installed into a single VM. In this case, Fedora was chosen over CentOS or Ubuntu for the existence of a community-written installation guide ([24]) and scripting utilities to ease the process, even though the official *OpenStack Folsom* configuration guide [25] is written with Ubuntu in mind so many commands are not applicable to Fedora.

Both Essex and Folsom version were tested. The reason was that in spite of being Essex the officially supported version in Fedora 17, the quick installation guide suggested using the latest version available — Folsom by December 2012 — enabling a testing repository to that end. The degree in maturity observed moving from Folsom to Essex was startling: not only the web interface had been revamped giving it a more thorough look that also reflected much better the underlying state of the infrastructure, under the hood, the core module had also been split in smaller functional pieces that could be easily distributed across the cluster. In Essex, when dealing with creating instances in the cloud, if there were a problem while the networking interfaces were being brought up, making it impossible for them to obtain IP addresses, the web interface would hang in a state that would not allow to destroy the instance requiring the invocation of console commands. This hanging problem is solved in Folsom.

An execution environment for both versions was spawned following the next list of actions (see figure 3.4):

- A single VM was created to hold OpenStack: 1 GB of RAM, 10 GB for the HDD, APIC and ACPI.
- Fedora 17 was downloaded, installed and yum-updated as usual.
- The VM was named *openstack*.
- The *openstack* user was registered as well as added to the sudoers list.
- `acpid` package was installed.
- *SELinux* was stopped.
- A full clone of the VM — which includes the virtual drive — was made to test both versions.
- The quick installation guide as well as the official one, omitting Swift, were followed to complete the process.

3.3 Veredict

What follows is the listing explaining the findings in the comparison between the frameworks according to the methodology explained in the beginning of this section.

Installation: Without a doubt OpenNebula stands out. The installation guide is the lighter and shortest to follow with difference. Carries, however, the inherent problem of hiding what is going on when it is being configured for the first time, potentially hardening the resolution of issues that might appear in the future.

Configuration and Management: Growing the supporting cluster requires, on every cloud tested, that a compatible hypervisor be installed on any node added. CloudStack and OpenNebula offer a more transparent management interface to better keep in check the physical infrastructure. OpenStack displays the most limited web interface.

Hypervisor: Regarding hypervisor support, OpenStack clearly surpasses both CloudStack and OpenNebula. Yet, they all support the most widely used hypervisors — KVM, Xen, Xen variants and VMware and variants — in production deployments.

Storage: The three of them support a broad assortment of data back-end controllers. But, in this case, it is important to highlight the effort OpenStack is involved in to introduce Amazon S3 compatibility in its deployments. Swift is the OpenStack component granting fault tolerance and high availability storage, mimicking Amazon S3, relying on data replication and balancing among other techniques. CloudStack advanced installation guide [21] describes a first approach toward configuring Swift as secondary storage for the cloud. This fact speaks volumes about the maturity and importance of Swift, an OpenStack module.

Documentation: None of the three can boast about exhaustive official installation guides. Every framework has had its own exposure to different linux distributions, so the coverage they offer of them varies to the point of mistaking module names, e.g. both CloudStack manuals are more easily followed using CentOS as base operating system and XenServer as hypervisor. OpenStack provides installation manuals supporting both Red Hat and Debian derivatives, but for Fedora the name of the documented backing services does not correspond to the real ones; not such thing happened for Ubuntu. Nonetheless, inaccuracies are trivial to cope with and the manuals are deep enough to deploy in production.

Community: Even though it may seem unimportant, the community is vital for developing and supporting the frameworks. They are, at the very least, partially open-sourced, so a lively community translates into higher usage rates, more rigorous documentation, more bugs squashed, etc. While it is hard to assess the magnitude of an online community

from the outside, it is interesting to highlight the nourishing that OpenStack is continuously receiving from Red Hat and Canonical: there is no technical keynote or conference in which OpenStack is not appointed.

3.3.1 OpenStack Folsom

The IaaS Cloud that has been chosen is OpenStack Folsom. The lengthy installation guides, the community support, the backing by two large software companies, the real deployments in production (from HP, Dell, Intel, Rackspace, etc.), the modular configuration, the completeness of the implementation (OCCI APIs, S3, EC2, Swift, etc.) and the *official* support to deploying the cloud over a virtual cluster for testing have unbalanced the comparison in its favor.

Chapter 4

OpenStack Folsom

The current section intends to detail the IaaS Cloud implementation that has been chosen: OpenStack. Initially, a global vision will be given to the reader, to progressively focus on its constituents modules' responsibilities and how they collaborate to maintain the service running.

4.1 Global Architecture

Figure 4.1 shows the three basic operational components of OpenStack Folsom:

Functional Core: OpenStack Compute, OpenStack Quantum and OpenStack Storage (Cinder and Swift).

Web Management Interface: OpenStack Horizon.

Shared Services: OpenStack Glance, OpenStack Keystone and other related services like a DBMS for persisting meta-data or a messaging queue.

The different components have been devised in a shared nothing fashion. This provides the cloud admin the flexibility required to distribute the modules over the cluster as pleased. An example of a particular OpenStack

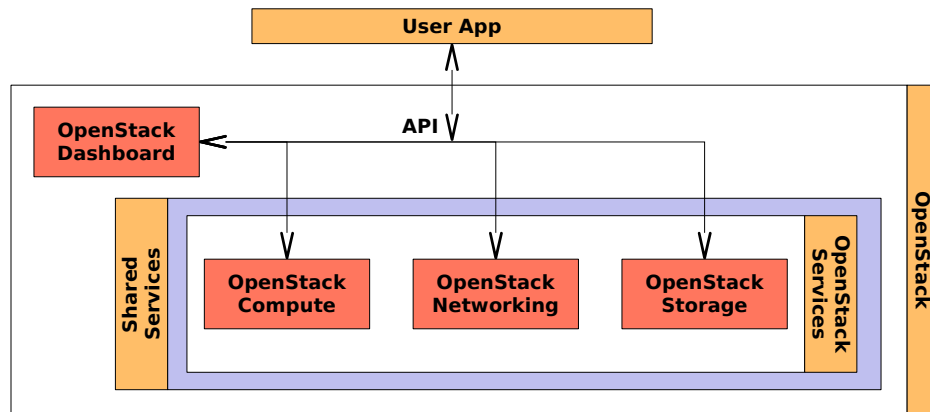


Figure 4.1: OpenStack Architecture

deployment is shown in figure 4.2; OpenStack’s own modules are displayed in red, supporting services are shown in violet. What it is missing from the diagram, for clarity, is the asynchronous queue that mediates inter-module communication. Qpid and RabbitMQ are the two queue implementations that are officially documented, being the former the one that we used in our test deployment.

4.2 Horizon

Horizon represents the fundamental window to set up the cloud. As discussed in the previous section, Horizon does not currently — as of Folsom version — present a global view of the physical infrastructure, leaving the user in the dark in this respect. Horizon is written in Python on top of Django, the web framework. Django itself relies on a web server like `httpd` to expose static files, uses a caching mechanism (`memcached`) to speedup load times and a terminal embedding (`noVNC`) system to view the output of the virtual graphic card directly on Horizon.

To manage and create instances in the cloud, OpenStack gives the cloud admin the ability to register authorization roles that will let the users con-

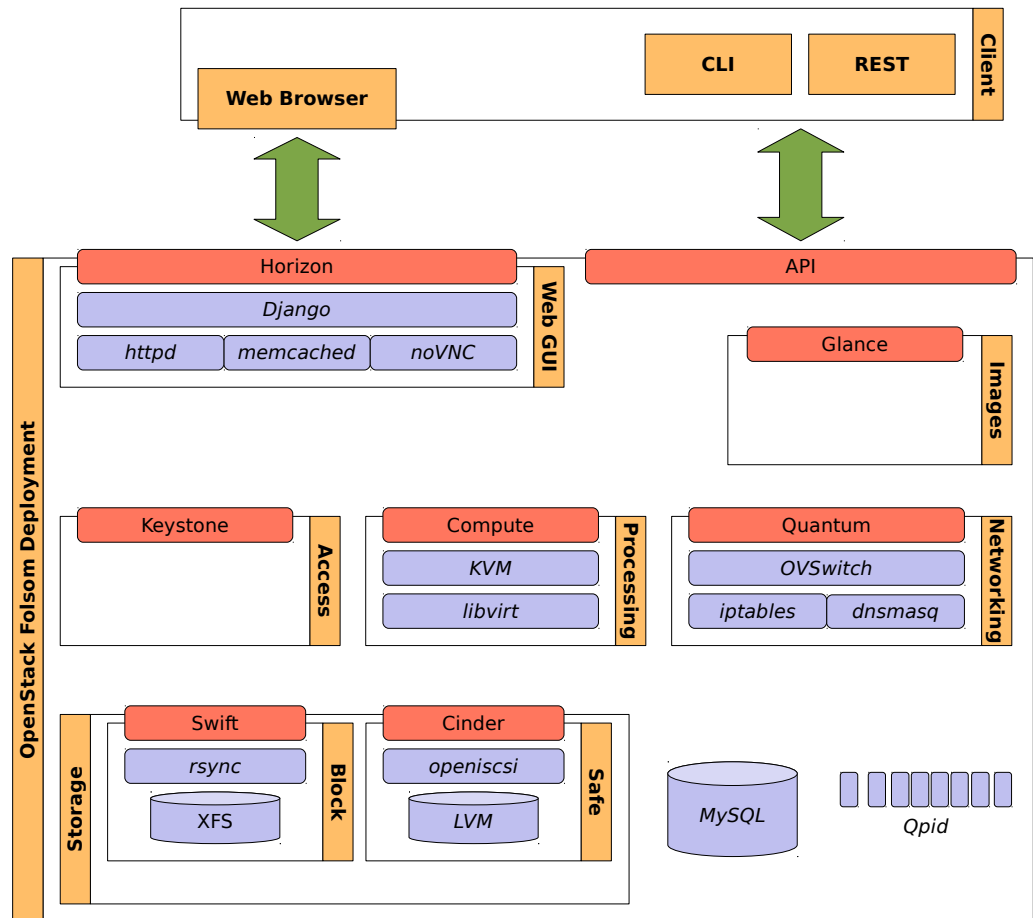


Figure 4.2: Example of an OpenStack Folsom deployment

sume those services whose role give access to. While the admin is allowed to sign up custom roles, two roles that ship the distribution are the *Cloud Admin* and the *Cloud Member*.

A user granted the admin role will be able to manage:

Tenants: Create, delete, member users, alter quotas, etc.

Users: Create, modify or delete.

OS Images: List, remove or modify meta-data.

Instances: Reset, shutdown, suspend, print log on screen, etc.

Volumes: Create, list, attach to an instance, etc.

Networks: Create, modify or delete.

A user granted the member role will be able to:

Status: Quota, resources, etc.

Instances: create, shutdown, reset, suspend, print log, create image from a running instance (snapshot), etc.

Volumes: List, create, modify, attach to an instance, create a volume snapshot, etc.

Images: Create, list, delete, modify, etc.

Networking: Manage public IPs (floating IPs).

Security groups: Create, delete or modify security rules.

Keypairs: Create, modify or delete.

4.3 Keystone

Keystone is the central security check point and information repository storing information needed to access the cloud installed services. It verifies, before each request, user credentials and authorizations in OpenStack services. Keystone divides this functionality in two parts: on the one hand user control, on the other service catalog.

To deal with users, Keystone assigns them tenants or projects. Users, as discussed above, are granted the membership to a tenant and a service quota they will have to adhere to; they are also restricted to the tenant quota.

To organize the catalog at hand Keystone defines two other concepts within the service catalog: *Services* and *endpoints*. A service in the catalog is a mere abstract description of an exploitable cloud feature by the user. The particular implementation of the service is managed by the set of endpoints associated to it. Said collection contains every piece of information that is required for users to consume the services. Figure 4.3 shows a Sequence Diagram portraying the interchanged messages between the different entities taking part to consume a service: *Create a new instance*.

Stemming from the fact the Horizon exposes only a part of OpenStack functionality, to help dealing with security Keystone installs a CLI tool to interact with the REST service in charge of administrative operations. Issuing certain commands to Keystone through a terminal requires the knowledge of the admin token, which has to be conveniently secured, or the login credentials of a user with the admin role. Lastly, it should be noted that Keystone uses a data base to store user access credentials and the service catalog metadata.

4.4 Quantum

Starting with Folsom, Quantum is the module to manage virtual networking. It was introduced to separate the networking part from the computing part, held together in **Compute** module. Certainly, the fact that it had been

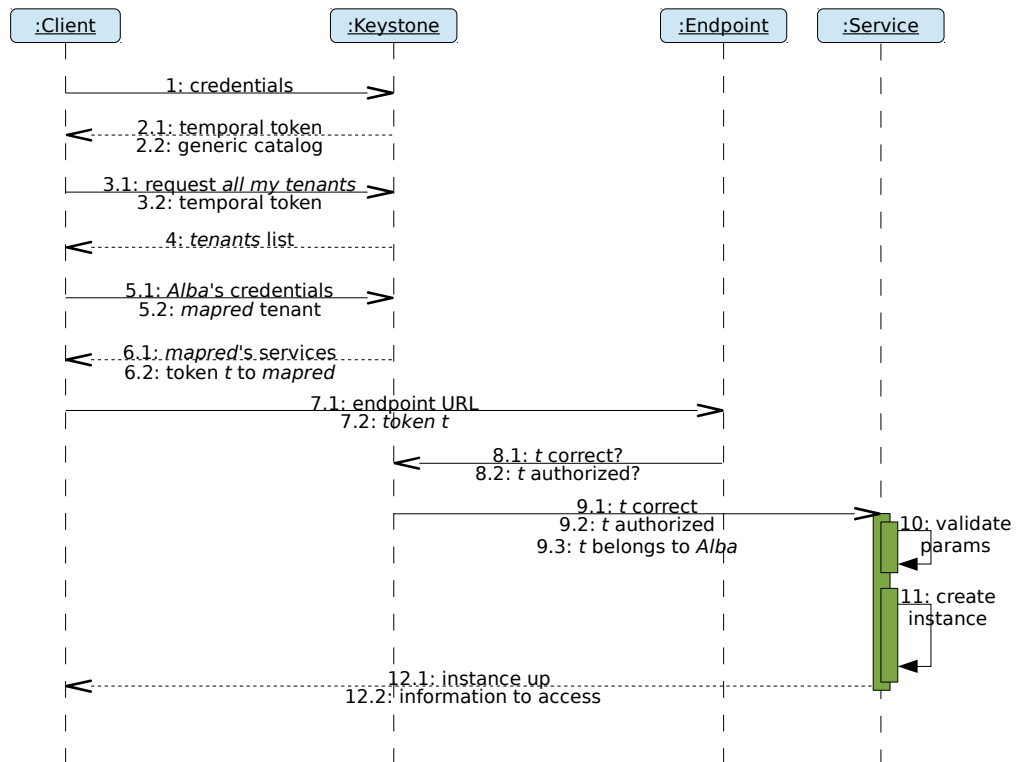


Figure 4.3: Sequence Diagram — create instance

refactored out demonstrates OpenStack's evolving model toward a more coherent less coupled functional allocation; and as it is independent, it could be configured in a dedicated node.

To bring virtual networking into existence Quantum banks on external plug-ins. Two of those plug-ins whose usage is covered in the official Quantum administrator manual ([26]) are `OpenVSwitch` and `LinuxBridge`. Additionally, Quantum relies on `iptables` to configure routing rules and firewall, `dnsmasq` for the *DNS*, the *DHCP* and the *NAT*.

Figure 4.4 pictures a topology example of a virtual network. On it, $30.0.0.X$ represent public IPs and $10.0.X.Y$ private. This virtual network assigns a virtual router to each tenant but more could be added with ease. Private IP overlapping over different networks is possible as expected ($10.0.0.2$). The routers public IPs — they could be assigned more external interfaces — must be taken from the external network ($30.0.0.0$).

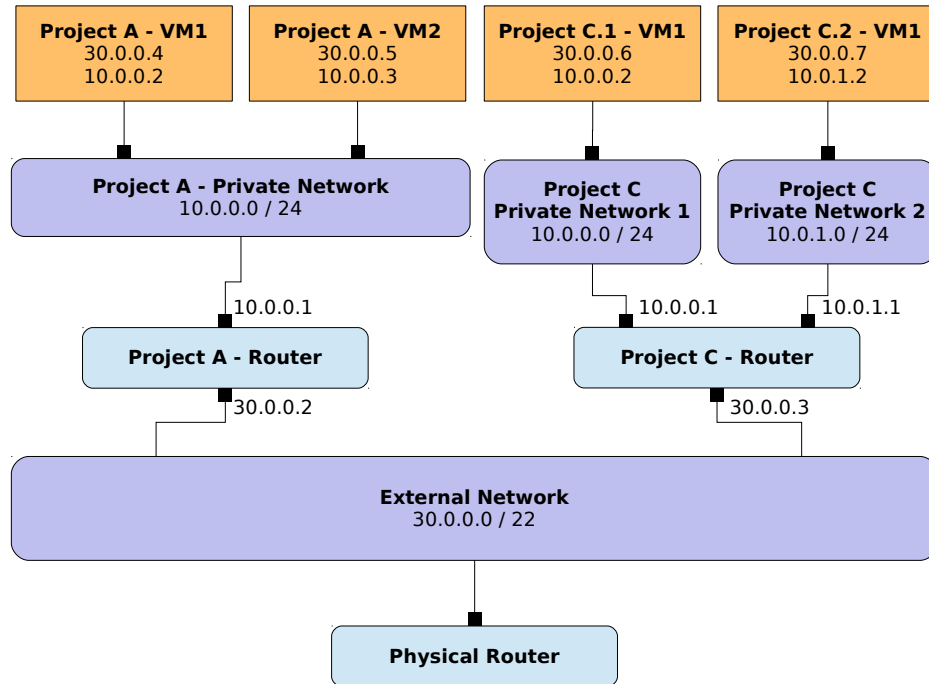


Figure 4.4: Virtual network deployment with Quantum

4.5 Compute

Compute is the central module. Its duty entails orchestrating the global workings in the cloud, delegating each particular function to the service on charge. In the end, Compute will let a logged user start virtual instances, which will draw their VCPU, VRAM and VHDD from the physical cluster. Yet required, Compute does not contain a virtualization package. The approach is to delegate infrastructure provision to a hypervisor found typically, but not restricted to, in the same node. To expose this on-demand computational service, Compute implements a REST API so that users can control their instances' life cycle directly from a REST client (like the CLI tools that accompany Compute).

To create an instance in effect, Compute will communicate with other modules within the cloud to orchestrate the execution and, finally, it will pass the request to the most suitable cluster node's hypervisor — most suitable according to the cloud-defined rule set — that will bring up the VM. Some of the supporting services are described below.

Keystone: Collates credentials y authorizes requests.

Glance: Selects the OS image that will be used to start the VM.

Quantum: Grants private and public IPs as well as manages instance network traffic.

Cinder: Manages block storage and on-line volume attachments.

Qpid: Handles message interchange between Keystone, Quantum, Glance and/or Cinder.

As it has been discussed all along, if there is some trait that aligns different IaaS Cloud implementations is their flexibility. Users' computational needs are as diverse as they are changing and therefore they expect to be given the chance to define virtual infrastructure adapting to those needs. In OpenStack, each possible particular configuration instance will take its

VCPU, RAM and VHDD from a cluster host, and the users will be allowed to shape those to their requirements with ease.

4.6 Glance

Glance is OpenStack's OS image storage service. Glance may be configured to drive images stored in a myriad of backends, ranging from Swift to an HTTP-addressable location. As happens with every other OpenStack module, Glance relies on Keystone to grant access to the images, and coordinates its operation with Compute to put them in execution on demand.

Glance supports a good number of image and container types — this fact being merely informative to the Cloud framework as it is the hypervisor who would have to support the particular combination image type, container type —, and they are stored as metadata linked to the image in Glance.

4.7 Storage

OpenStack provides three main options regarding storage types:

Ephemeral: The size of the drive hosting the root file system is set following the particular flavor parameters when the VM starts. The files contained in this file system are those present in the image file stored in Glance. Any alteration to this file system will only persist the execution of the VM. Any change on the image files is written temporarily to be discarded as soon as the VM is shut down.

Block: By making use of storage volumes managed by Cinder with *LVM* (*Logical Volume Manager*) OpenStack provides the ability to attach indeterminably-sized logical volumes to instances on-demand. This store kind guaranties that information is preserved between VM executions. However, this method carries an important handicap, that of being unable to attach a single volume to two different instances at

the time. High availability or data safety on failure are not supported, as data is stored in a single place. A backup or RAID policies may be established to get over these limitations but they are discouraged as OpenStack has its own module to deal with them.

Safe: Swift manages a safe distributed storage banking on controlled replication allowing for high availability deployments that overcome hard drive's inherent fragility. Swift draws on `rsync` to synchronize *XFS* partitions.

4.7.1 Cinder

Cinder is the OpenStack module that takes care of virtual block storage devices — functionally similar to Amazon's *EBS* (*Elastic Block Storage*). Cinder uses an *iSCSI* implementation (`open-iscsi`) and LVM to manage operations on the volumes. Creation, attachment and detachment, and logical volume removal is directly controlled through Horizon.

Those persistent virtual blocks are administered as logical volumes pertaining to a volume group controlled by Cinder. Cinder, though, shall not be used to create a shared medium to instances, as *NFS* (*Network File System*) or a *SAN* (*Storage Area Network*) solution do; for a single volume cannot be attached to different instances at the same time. An interesting option that Cinder opens is using a logical volume to boot instances, therefore sharing the set of files contained in the image among them.

4.7.2 Swift

Just as happened with Cinder, Swift cannot be framed into traditional shared networking nor be compared to Cinder: Swift covers a different functional demand. Swift is defined as “a scalable object storage system where logged users control their store buckets uploading, downloading or deleting files to their will” [27]. Swift may be conceived as a functional clone to Amazon's S3 or Eucalyptus' Walrus, implementing a partially

Amazon-compatible REST API. Central to Swift's implementation is replication.

Replication

Scalability, fault tolerance, high availability, safety, storage and load balancing are some distinguishing features of Swift's. As discussed, high availability and fault tolerance are implemented with replication. Replication is a mechanism by which a distributed system keeps block copies at different locations of the deployment to guarantee better performance and limit failure impact.

Within Swift, replication processes on every *Object Server* — any node in the cluster configured to support them — periodically compare their local blocks with remote replicas to collate their update state. Comparing replicas states is as costly a process as it is often, thus *Hash lists* and *watermarks* are used to improve comparison time. Replication is transparent to the user and Rsync or HTTP transport replica payload across the cluster. When a new Swift node is added to the cluster, replica distribution becomes unbalanced and will trigger automatic rebalancing. When it be synchronized with the cluster, the new node will be able to respond to data requests.

Updaters y Auditors

Other supporting services that complete the functional circle defined for Swift are *Updaters* and *Auditors*. The former act when a replica synchronization error is raised or when Object Servers load is so high to make a data request stay unreplied. It happens then that the execution of this operation is delayed and queued, being serviced by the Updater process at a later time. Auditors continually scan the file system looking for integrity failures in objects or buckets. If an inconsistency were to be found, the incoherent entity would quarantined and replicas would be made anew.

Chapter 5

Hadoop

This chapter tries to expose with simplicity the defining fundamentals of Hadoop architecture. Initially abstract concepts will be introduced to give way to more particular and deep ideas that explore Hadoop implementation of the MapReduce model in two layers: the processing and the storage subsystem.

5.1 The Beginnings

Hadoop roots its origins in *Apache Nutch*, Mike Cafarella and Doug Cutting's implementation of an open source web index and search engine. Nutch project began in 2002. In spite of the Internet being notoriously smaller at the time, Nutch's underlying technology was unable to make it scale to manage the billion pages that comprised the *old* Internet. But in 2003 Google publishes a research paper introducing *GFS* (*Google File System*) [28], a file system to be used across their clusters of commodity pcs that greatly simplified its deployment. Nutch will inherit a large part of the concepts detailed there translated in their own distributed file system implementation (*NDFS*).

Also in 2004 appeared another publication [2] that presented MapReduce, bringing about successive efforts to port Nutch algorithms to adapt to the emerging model. In mid 2005 most of Nutch code run following MapReduce

guidelines over NDFS.

Both NDFS and Nutch MapReduce implementation were generic enough to be used without refactoring beyond web page indexing. In 2006 an unrelated project was constituted to extend Nutch's potentially reusable parts to widen their applicability context. This project was called Hadoop. In 2008 *Yahoo!* announced that the index for their search engine in production was being continually refreshed by 10,000 Hadoop nodes. This same year Hadoop is brought out to the world becoming an Apache-backed project corroborating its success.

Nowadays, Hadoop is without doubt the MapReduce implementation most widely used by a broad range of companies.

5.2 General Hadoop Architecture

Hadoop composition differentiates four modules:

Hadoop Common: A module containing the parts used across the implementation. It is mainly comprised of scripts and configuration tools.

Hadoop MapReduce: The module implementing the MapReduce processing model.

Hadoop YARN: A general purpose framework abstracted from Hadoop MapReduce. It is employed to manage resources and schedule executions in distributed environments.

Hadoop DFS: The distributed file system sustaining inputs and outputs from Hadoop clusters.

Hadoop architecture corresponds to the *Master-Worker* archetype where two roles on each cluster appear: a unique Master and various Workers. These roles, and thus responsibilities, are fixed to different nodes by the cluster admin. If necessary, e.g. for maintenance, the admin may freely reset the roles to new cluster nodes, only requiring job resubmissions if the Master role were reassigned.

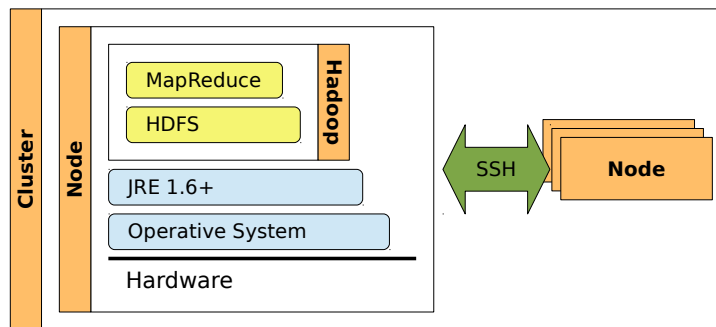


Figure 5.1: Hadoop over HDFS

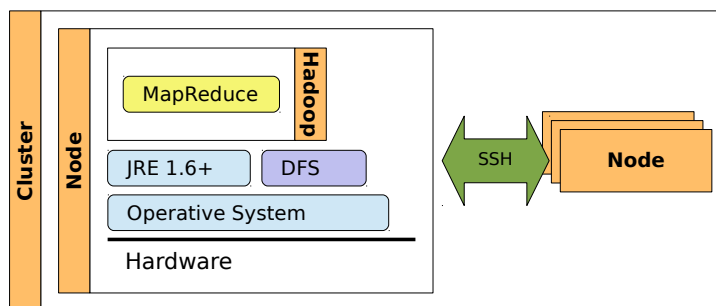


Figure 5.2: Hadoop over another DFS

This section will almost exclusively center around MapReduce and Hadoop DFS (HDFS) modules to expose the functionality covered with Hadoop. Hadoop YARN, as discussed, is a subsystem resulting from the isolation of scheduling and processing, both found together in the old — pre Hadoop 2 — MapReduce module, retaining task distribution and planning within YARN. This way, YARN is allowed to untie from Hadoop allowing for deployments where YARN orchestrates an implementation-agnostic working set. As of this writing, Hadoop YARN is still an alpha version.

Figures 5.1 and 5.2 exhibit a high level vision of Hadoop architecture. Figure 5.1 shows an hypothetical deployment with HDFS. Figure 5.2 shows a particular Hadoop installation with another supporting distributed file system.

From the figures it can be deduced that Hadoop runs atop a *Java Virtual Machine* (emphJVM), that MapReduce requires a DFS implementation to rely on and that inter-node communication is conveyed through *SSH* tunnels over TCP. Every module includes a web server (*Jetty*) to ease collecting and reporting status information.

5.3 Hadoop Distributed File System

HDFS has been designed to act as archival repository for huge masses of data whose main access pattern be *write-one read-many*. While it is no requisite for a data query to follow this access pattern, HDFS performance shines with read queries in batch mode. The underlying infrastructure, again, is composed by commodity pcs that HDFS manages to transform into a reliable, scalable, fault tolerant, self-balancing and network traffic reducing data store. Yet, as individual nodes are regular pcs, in HDFS converge some operating limitations:

- High data access time. HDFS prioritizes large reads in batches and thus, reading small files is generally discouraged. As complement though, HDFS delivers *very* high throughput by leveraging parallel reads on the cluster.
- High data writing time on many small files. Writing a file changes a block. Every file that be smaller than the HDFS block size must be persisted in a single block. That changed block must be send out across the network to keep the file system consistently updated. Thus, appending to many files requires updates in many blocks that would require synchronizing over the network.
- Multiple writes to the same file or “*not-append*” operations are not supported. HDFS is not a *POSIX*-compilant file system implementing only a set of operations to try to maximize data throughput in distributed environments.

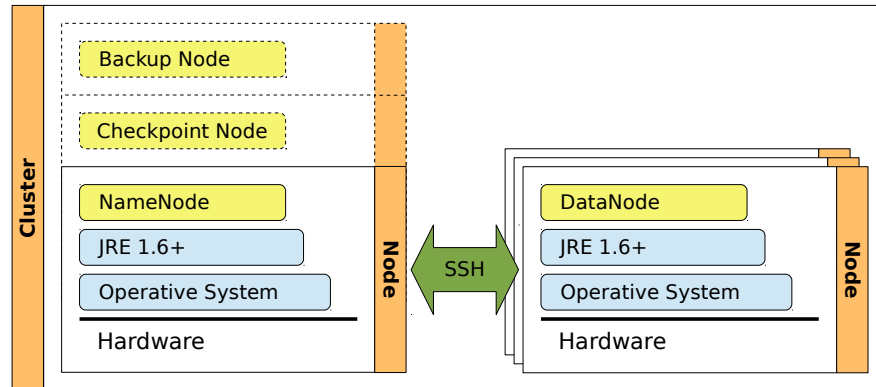


Figure 5.3: Typical HDFS deployment

To organize storage, HDFS takes from traditional operating systems the concept of block: an abstraction of the particular drive structure with a double purpose:

Lowering DFS complexity: Writing a block comprises storing data and meta-data, and handling information related to locate those data on disk. Using the block as the minimum organizational structure simplifies location expressions.

Incrementing flexibility: files are free to grow over the size of an HDFS block.

To lay local drive blocks in position, HDFS makes use of two processes: the *DataNode* and the *NameNode*. Besides, as support, from Hadoop 0.21.0 onward it is permitted the optional deploying of a *Backup Node* and a *Checkpoint Node* in the same cluster.

5.3.1 Node Roles

Figure 5.3 shows a layered down HDFS deployment. In dotted line appear both the *Backup Node* and the *Checkpoint Node*.

DataNode

DataNodes are those processes within nodes that handle the storage of HDFS blocks in their local drives. Every time a write to a file in a block succeeds, the DataNode in charge of the operation signals its supervising NameNode so that it could keep track the modifications to the DFS as they happen.

NameNode

The NameNode is the process appointed to deal with the name space of the cluster. It has to handle the file system tree and meta-data that make possible recovering data from the DFS. It is such an important process that if it went down, every piece of data in HDFS would get lost, rendering impossible to match files with their container blocks. Therefore, the NameNode is seldom deployed with no Checkpoint Nodes or Backup Nodes in case the HDFS data were not stored elsewhere.

The information about the file system, meta-data, is persisted to the NameNode both in memory and disk. In this latter form, meta-data is managed in two files: one contains the name space of the file system as an image (**fsimage**), the other progressively appends the changes to the fsimage (*edits*) as a log. When the DataNode starts, it compiles an fsimage afresh by merging the existing fsimage with the edits file. As soon as changes to the HDFS are reported from the NameNodes, the DataNode will update the edits log but without touching the fsimage. In a typical secured deployment, the edits file is kept in memory within the NameNode host computer, on the local file system and remotely via NFS.

Checkpoint Node

The Checkpoint Node is the means by which failure in the NameNode poses a smaller threat to the integrity of the data within HDFS. The idea behind the Checkpoint Node is to maintain a separate copy of the edits file so that it could periodically merge the fsimage with the edits, generating an updated

fsimage to be uploaded to the NameNode. It should be noted that the Checkpoint Node does not listen to the network to record the modifications to the DFS in the edits file, it fetches the copy held by the NameNode itself. When the merging process be finished, the NameNode will be submitted the newly created fsimage, clearing out the old copy and resetting the edits file.

Backup Node

The Backup Node provides the same functionality as the Checkpoint Node by periodically creating restoration points of the name space but through a different approach. In this case, the Backup Node will download the fsimage file from the *backed up* NameNode when booted, just like the Checkpoint Node, but it will receive notifications from the DataNodes on each modification to the HDFS, and will manage itself the merging between the fsimage and the edits file effectively creating a new fsimage, the same behavior as the NameNode.

Compared to the Checkpoint Node this process draws less bandwidth due to the fact that it does not require to download the fsimage nor the edits from the NameNode in order to stay synchronized.

The Hadoop version used in the VM of this project — version 1.0.4 — allows only a Backup Node per NameNode or multiple Checkpoint Nodes, not both. It should be noted that including a Backup Node in a cluster gives the possibility to run the NameNode with no persistent storage, i.e. with no allocated space in the local drive to write the fsimage and the edits, making the Backup Node the sole responsible for the duty.

5.3.2 Network Topology

One of the fundamental parts to a file system in a distributed environment is the capacity to provide a transparent mechanism by which information is persisted securely while, at the same, a certain amount of performance is maintained. HDFS makes use of an already discussed technique, *replication*, that provides high performance, scalability and fault tolerance. Besides,

HDFS has been devised to reduce the network congestion that stems from regular operation giving special emphasis to the way in which data is to be distributed in the cluster, controlling replica count or distance between replicas, among other variables, to concrete the block allocation policy.

Node Distance

To support the features that have been mentioned, it is fundamental for the NameNode to possess some information on the participating DataNodes physical location. With that information, the NameNode would be able to balance the *physical distance* of the replicas, realizing that, generally, fault tolerance increases with replica distance but the bandwidth available to transfer replicas narrows. Thus, the NameNode will have to distribute the replicas across the cluster with a procedure that make the cluster safe without clogging the network.

To effectively calculate the physical distance between two replicas, the NameNode will use an approximate measure of the physical distance of the nodes storing the replicas. Therefore, the problem lies in finding a formula to calculate the distance between any two nodes in the cluster.

IP-based networks arrange participating nodes in an abstract tree structure that is simplified in four octets. Without deepening any further on this topic as it is clearly out of scope, it should be noted that, in properly configured networks, the closer two nodes are physically the smaller the distance between their respective routing gateways. This idea may be used recursively until the *same* routing gateway was reached up from the initial nodes. So, the actual distance between two replicas is *the number of steps required to find a common routing gateway, starting from the two nodes that stored the replicas*. The example below will clarify the procedure.

The way in which the NameNode calculates the distance among replicas lacks an important feature that would render the cluster failure-prone. In order to keep HDFS fault tolerant, the NameNode has to deal with halting nodes, clogged networks, etc. If a stalled node were the gateway to a set of

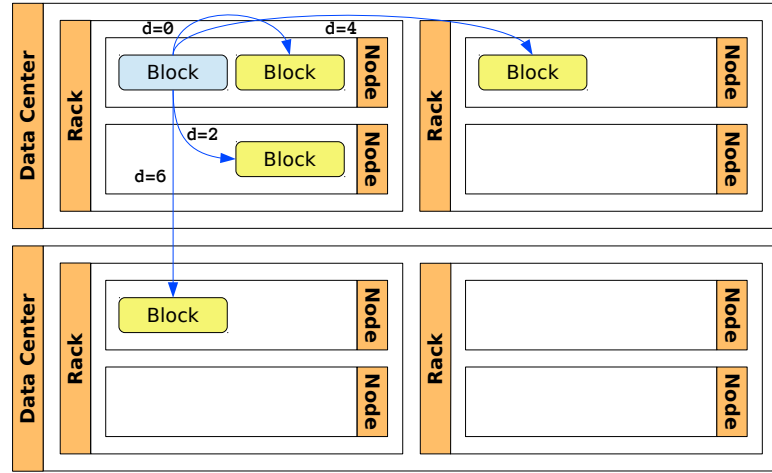


Figure 5.4: Replica distance

nodes, e.g the nodes in the same rack, then the access to any replica within that cluster would be impossible, and worse, if every replica of an individual block were stored in this rack, the block would be inaccessible for the duration of the outage in the gateway. Therefore, the distance alone cannot solve this distributing problems.

HDFS overcomes this limitations by mapping every IP address to a tuple with as many components as levels there were in the network — typically three-tuples (*data center*, *rack*, *node*). Figure 5.4 shows the most common values of the replica distance. With the help of the figure the procedure to get to $d = 4$ will be succinctly explained.

The blue node holds the original block and the yellow one — pointed by the $d = 4$ -labeled edge — a replica. Both blocks are stored in nodes within the same data center but in different racks, thus, their nearest common ancestor gateway would be the one managing the routing between those two racks. Besides, every rack has an internal gateway that must be traversed in order for the internal traffic to exit the rack. So, each node would have to take 2 steps to get to their common gateway — from node to *in-rack* gateway to *off-rack* gateway — coming to 4 steps.

Replication

Replication is a NameNode-controlled technique that relies on knowing the network topology to structure hierarchically the nodes in the cluster. Together with replica distance the NameNode is capable of enforcing a policy that balances data safety and throughput. What follows is the description of the method to distribute replicas over the cluster.

When a new block is written to HDFS — or when an existing one is appended new data — a configured number of replicas have to be made and distributed. By default, the NameNode will maintain updated three replicas of each block. The first block is placed in a node at random prioritizing those not too full nor too busy. Then, the first replica is stored in the same node as the original block. The second replica is sent off-rack to a node at random, and finally, the third replica is stored in the same rack as second but in a different node. If the *replica factor* — the total number of replicas every block will have — were made higher, the successive replicas would be placed, again, in the same rack as the second one but in different nodes whenever possible.

This method provides the desirable equilibrium between fault tolerance (by assuring two copies of a block exist in different racks), consumed bandwidth (writing a new block would only need replicas to go through a single gateway as subsequent replicas are made within the rack), read performance (making it possible for every read request to be fulfilled from nodes in two racks) and balanced data distribution by executing the process just described. Figure 5.5 exemplifies a factor 3 replication. The number in brackets indicates the order in which the replicas where created. The original block is the blue-colored one.

5.4 Hadoop MapReduce

In a typical Hadoop deployment, the execution module is lied over HDFS. This execution module is basically an implementation of the ideas exposed

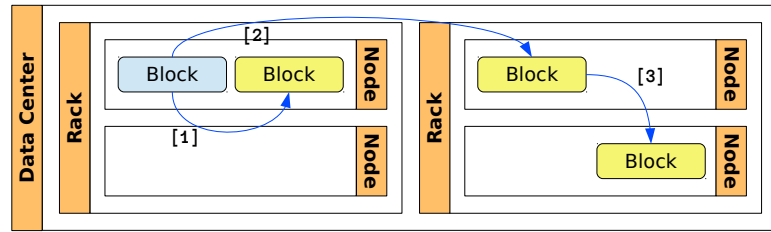


Figure 5.5: Replication factor 3

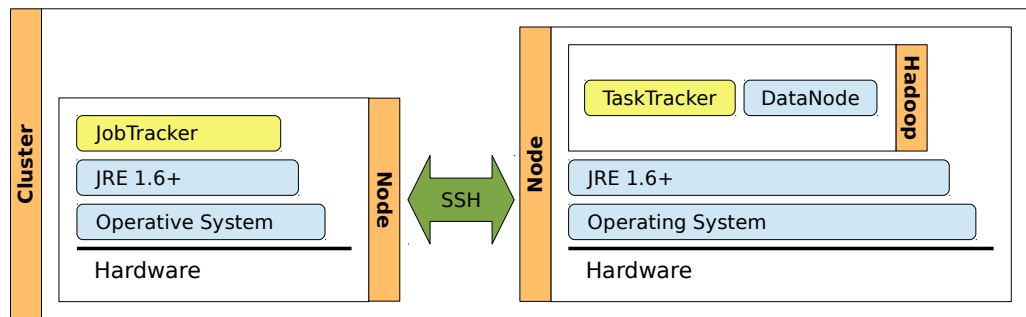


Figure 5.6: Hadoop MapReduce with HDFS

in the Google paper on MapReduce as a new programming paradigm [2]. The core ideas that have been discussed for HDFS hold in this section, but with a different flavor. Hadoop MapReduce architecture is also based on the Master-Worker model and so appear two roles: a Master and a Worker.

The Master role is functionally covered by a process called the *JobTracker*, whose main responsibility is task scheduling to workers. Complementarily, the worker role is implemented in the *TaskTracker* which will actually execute the individual tasks reporting their progress back to the JobTracker.

5.4.1 Node Roles

Figure 5.6 shows the JobTracker and the TaskTracker in context, using HDFS and file system. To complete the picture HDFS processes, the NameNode and the DataNode, would have to be included.

Just as with HDFS, the JobTracker and the TaskTracker will be subsequently dissected.

JobTracker

The JobTracker behaves in a very similar manner to the NameNode, but in this case the former handles the scheduling part of job execution. A regular workflow starts with a client submitting both the algorithm and data to the cluster. The client is returned an identification for the new job and it is its duty to upload the data to HDFS. When data be present in the cluster, the client would transfer the job configuration to Hadoop. The Jobtracker, when it had determined that the job should start, would initialize the job sending out to the TaskTrackers the code required to complete the execution. The TaskTrackers would be sent a subset of the work to be completed for the job; a task in MapReduce nomenclature.

Distributing the work among TaskTrackers is done following the principle of maximum locality, i.e. making TaskTrackers process the parts of the job that refer to data stored closer to them. Enforcing this policy would limit the amount of traffic flowing over the network and, at the same time, would shorten data access time.

Running in an environment prone to error, the JobTracker must define a method to deal with a myriad of failures with limited impact on job throughput. To that end, the JobTracker maintains a list detailing task status. This list is updated by the JobTracker process continually polling the TaskTrackers for new information on the tasks. If a TaskTracker were to stop responding after multiple tries, the JobTracker would mark the task failed to reschedule its execution at a later time.

More difficult to deal with would be the failure of the JobTracker process. The initial approach as proposed by Google in their so cited paper ([2]) is to discard unfinished jobs and try to reboot the service in an idle node, as only one JobTracker per cluster was supported. In Hadoop ecosystem *Zookeeper* provides the means for running multiple instances of JobTrackers within the

same cluster.

TaskTracker

The TaskTracker is the process that executes MapReduce algorithms on the data ideally stored in the node running the process. The TaskTracker will periodically communicate with the JobTracker to report status updates. As it has been discussed above, a failing TaskTracker will be unable to repeatedly *pong* the JobTracker's *ping*, effectively rendering the task inaccessible and due to rescheduling.

Chapter 6

A Private Cloud for MapReduce Applications

This chapter introduces a novel solution combining the virtual infrastructure managed with OpenStack and the Hadoop implementation of MapReduce to conform a powerful computational tandem. *qosh*, as this project has been called, will be described along this section moving from architecture to implementation.

6.1 Architecture

Figure 6.1 shows a high level portrait of *qosh* execution environment. The component that acts as interface between system and user is displayed on the right end. It abstracts the inherent difficulty in configuring the job execution context and in deploying the virtual cluster. Furthermore, *qosh* will keep track of submitted jobs, MapReduce `jar` files, input data and output results, with no need to walk the HDFS to search for data.

To streamline development, it is determined that the very first *qosh* version be tested on a personal computer. As shown on figure 6.1, Fedora Linux is installed atop the hardware and OpenStack Folsom set up within. *qosh* will draw on the infrastructure provided by the local OpenStack configura-

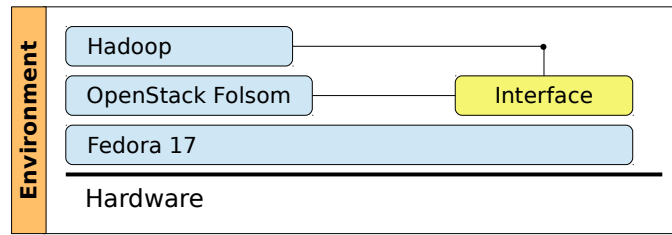


Figure 6.1: Global Architecture

tion to deploy virtual Hadoop clusters. While it would be better to make qosh more flexible allowing for remote infrastructure consumption, it would also become harder to test.

Figure 6.2 shows a more detailed view on qosh design details. The VM contains a Hadoop installation ready to be put to use as soon as it was started. The VM life cycle is managed by OpenStack and its execution environment shaped by KVM, the chosen hypervisor. Besides, HDFS has been used as temporal persistence layer while the results are not send back to the controller — the same machine in this testing environment. It shall be recalled that even though HDFS is a steady data store, the Hadoop VMs are created and removed for every work flow execution effectively destroying HDFS data at the end of processing each job. So, it is qosh’s job to orchestrate data extraction before shutting down the virtual cluster.

Figure 6.3 shows the modular decomposition of qosh orchestration module.

Compute: Acts as client to OpenStack REST API. It handles every interaction with the cloud decoupling qosh fromt the particular IaaS Cloud.

Django: Is used in qosh to let users manage their MapReduce executions with ease via a web interface.

Fabric: Is the Python library included in qosh to configure the virtual cluster deployment and destruction.

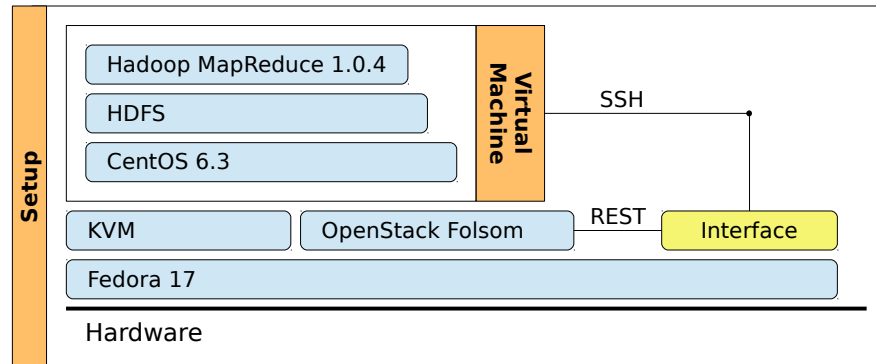


Figure 6.2: Detailed global architecture

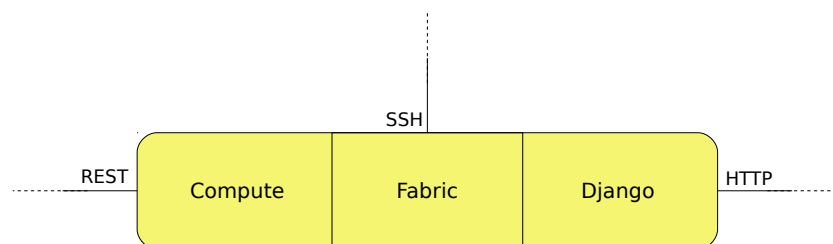


Figure 6.3: Core qosh modular decomposition

6.1.1 Design Diagrams

Below are shown the design diagrams that make up the section on high level overview of the project.

Django Components

Figure 6.4 portrays modules adjacent to Django to support its operation.

Apache httpd: Relied on to manage user interaction with OpenStack Dashboard, it may also be used to with qosh web interface. Initially, qosh depends on Python web server module to handle user requests, but *httpd* might be easily configured.

memcached: Is employed to cache web pages in order to speedup load times.

MySQL: Chosen relational DBMS to store job meta-data.

Use Cases Diagram

Figure 6.5 displays the set of use cases that have been considered for qosh. It reflects the five fundamental agents comprising the system.

Machine State Diagram

Figure 6.6 presents a summary on the navigation flow across the web interface. An example interaction is subsequently described.

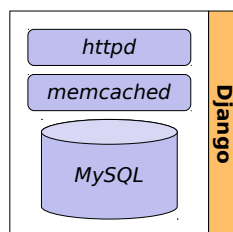


Figure 6.4: Django setup

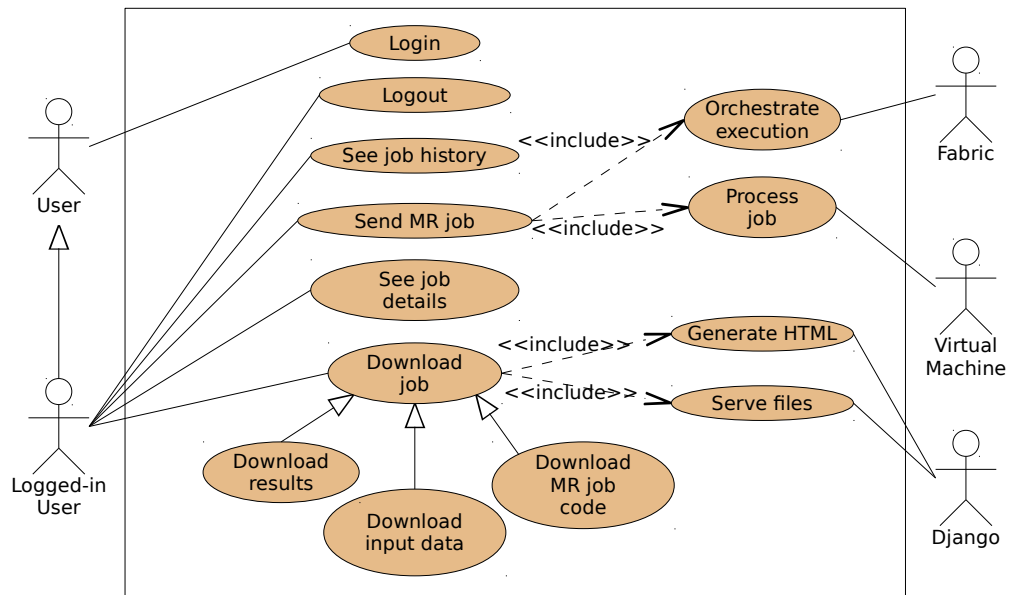


Figure 6.5: Use Cases Diagram

Initially, the user is presented the *Login* page so that he/she could log into qosh. If the supplied credentials were cleared — the user must be previously registered in Keystone as user/pass is shared with qosh —, the *Main* page is shown. From there the user may *Configure Job* or go over the *Job History* to get some *Job Details*.

Diagrama de Clases — Módulo Compute

La figura 6.7 muestra un pequeño Diagrama de Clases que describe las relaciones del cliente de acceso al API REST con algunos módulos de Fabric y Python.

json: *parsea* y maneja datos estructurados en formato JSON. El API REST de OpenStack entiende tanto XML como JSON.

Exception: clase Python que representa una excepción genérica.

ServiceError: extensión de **Exception**. Se utiliza como contenedor de los

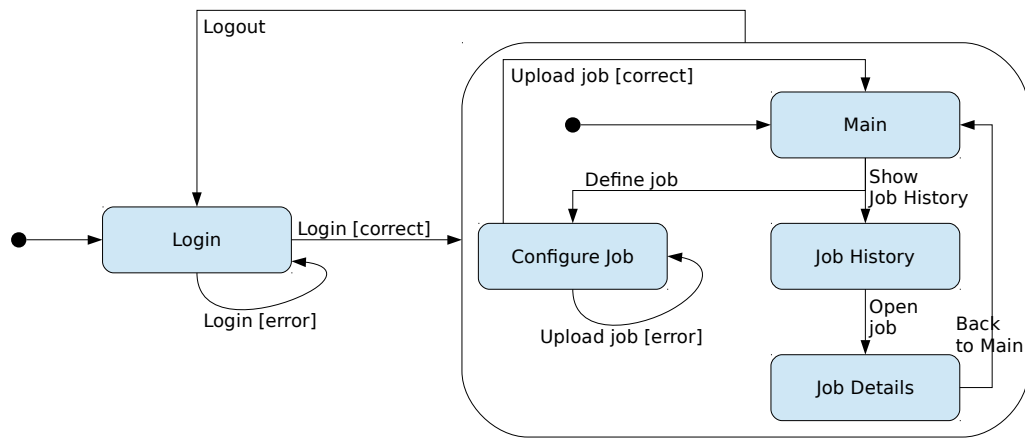


Figure 6.6: Web interface transitions

errores que pudiesen aparecer al interactuar con el servicio REST. Está compuesto por el código y la descripción del error HTTP que se produzca.

Environment: recoge las variables globales de configuración de la interfaz.

httplib: es el paquete Python que contiene las funciones y clases necesarias para establecer comunicaciones HTTP. Se usa en este proyecto para consumir el servicio REST.

6.2 Implementación

Los tres módulos que componen la interfaz se han escrito en Python. Las pruebas de la primera versión corren sobre las versiones 1.4.3 de Fabric, 1.4.2 de Django y 2.7.3 de Python. La configuración dentro de la máquina virtual, para permitir su acceso desde el cloud, se realiza con tres scripts, escritos en `bash-script`, que se activan en tres momentos diferentes del ciclo de vida del servidor virtual: antes de configurar la red, después de haber arrancado con éxito y antes de destruirse.

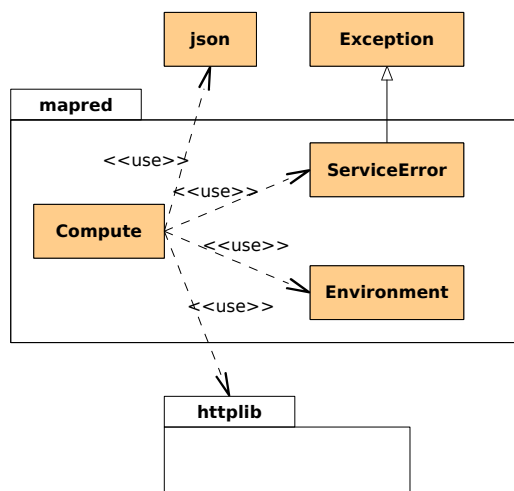


Figure 6.7: Diagrama de Clases del cliente de acceso REST (I)

6.2.1 Máquina virtual

Una parte importante del proyecto ha sido la construcción de la máquina virtual que contiene la instalación de Hadoop 1.0.4 y el JRE 1.7 de Oracle. Su configuración y puesta a punto no resultó un proceso muy complejo, pero sí lo suficientemente interesante y largo, por ser reutilizable para cualquier máquina virtual que se pretenda preparar para correr en un cloud, como para elaborar una lista con todos los pasos seguidos. El entorno de instalación es el citado MacBook 6,1 bajo un Fedora 17 actualizado.

- Se añadió con `yum` el *Virtual Machine Manager* (`virt-manager`) y sus librerías dependientes, entre las que destacamos: la librería de virtualización (`libvirt`), el módulo del kernel para el hypervisor (KVM) y un *wrapper* de ese módulo para gestionar su uso (QEMU).
- Utilizando el Virtual Machine Manager se configuró una máquina virtual con 1 GB de RAM, 4 GB de disco en formato *QCOW2* y tanto APIC como ACPI.
- Se prosiguió con la instalación en red de la última versión estable de

CentOS (la 6.3), eligiendo *Basic Server* como conjunto de paquetes instalados por defecto; en una sola partición ext4 y sin LVM.

- Al reiniciar, se pidió a yum que actualizase todos los paquetes del sistema.
- Se descargaron desde las webs oficiales las versiones comentadas del JRE y Hadoop y se instalaron con `rpm`.
- Se creó el usuario `hduser` con grupo principal `hadoop` y se retocaron los permisos de los ficheros relacionados con Hadoop —ficheros de configuración y scripts.
- Se configuró `sshd` para impedir la conexión como `root` y el acceso utilizando nombre de usuario y contraseña como credenciales; sólo se permite establecer túneles SSH usando la parte privada de la clave que inyecta OpenStack en la máquina virtual en cada arranque.
- Se escribieron tres scripts (`/etc/init.d/cloud-*`) para personalizar el comportamiento de la máquina virtual, como escribir en su lugar (`/home/hduser/.ssh/authorized_keys`) la clave pública de acceso SSH inyectada por el cloud.
- Se eliminaron, con `yum groups`, servicios no utilizados, como el *X-server*.
- Se compactó la imagen QCOW2 creando un fichero indeterminadamente grande, con `dd`, y comprimiendo, desde el anfitrión, con `qemu-img`.
- Con `qemu-nbd` se montó la imagen de la máquina virtual como un dispositivo de bloque en red y con `fdisk` se observaron el tamaño del bloque y el bloque inicial de la partición que contiene la instalación de Hadoop, para calcular su desplazamiento con respecto al inicio del disco virtual.

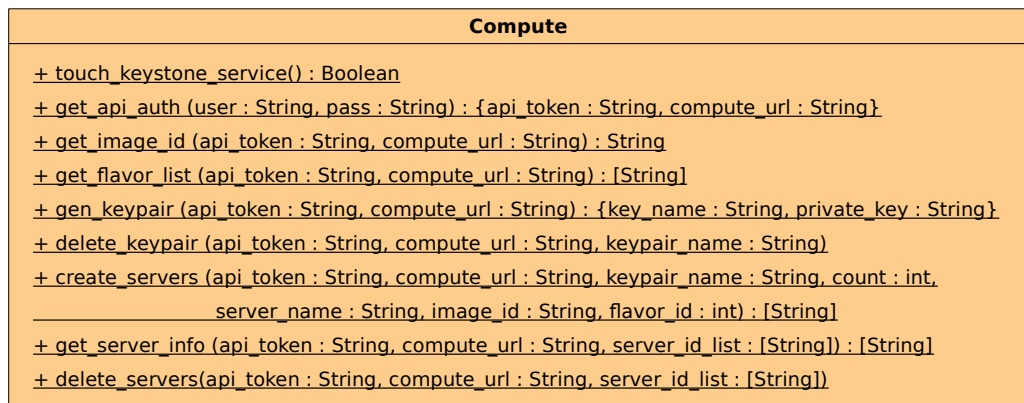


Figure 6.8: Diagrama de Clases del cliente de acceso REST (II)

- Finalmente, de nuevo con `qemu-nbd` se montó la imagen con el desplazamiento calculado, para extraer la única partición (la raíz) del disco virtual contenido en la imagen. Asimismo, se copiaron al anfitrión tanto el kernel como la *initram* de la instalación de CentOS.

6.2.2 Diagramas de implementación

Acto seguido se exponen aquellos diagramas que describen el proyecto desde un punto más cercano a su implementación.

Diagrama de Clases — Módulo Compute

La figura 6.8 expande el detalle del módulo `Compute` de la figura 6.7. El comportamiento de las funciones se deduce fácilmente de las firmas de la figura. Si cabe, destacar que la nomenclatura utilizada para describir algunos tipos de funciones —un tipo diccionario y un tipo lista, recuerda a la sintaxis de Python que declara esos tipos.

Lista: [TipoDeTodosLosValores]

Diccionario: {<Clave1> : <TipoValor1> , <Clave2> : <TipoValor2>}

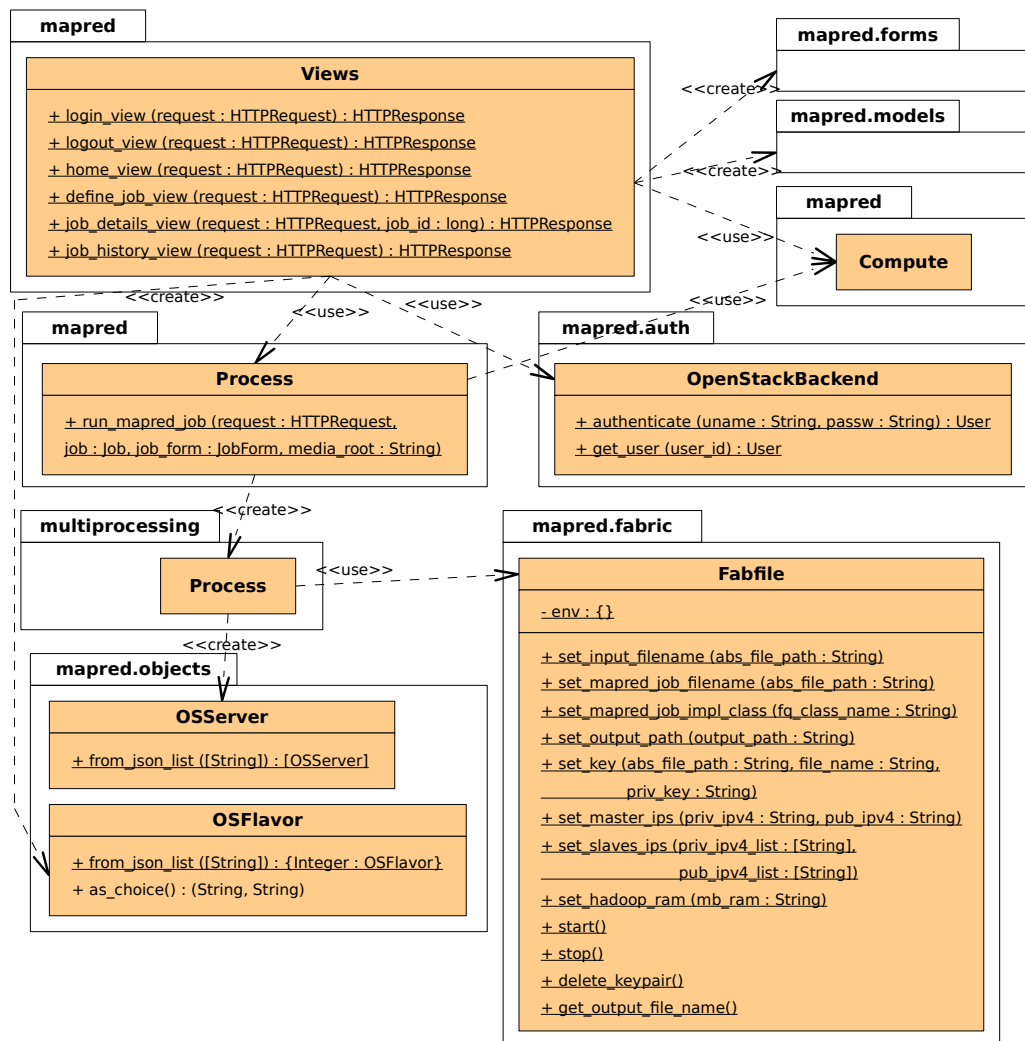


Figure 6.9: Diagrama de Clases — Django y Fabric

Diagrama de Clases — Django y Fabric

La figura 6.9 expone las relaciones entre los módulos más importantes de Django y Fabric. Veamos algunos comentarios sobre el diagrama.

Views: es una clase Django que implementa el comportamiento de cada *vista* (página web) de la interfaz del proyecto. Interactúa directamente con **Compute** para manejar la comunicación con el cloud.

Process: es una clase wrapper que encapsula tanto el comportamiento de creación de un proceso `multiprocessing.Process` como la función que ejecutará el proceso creado. Esta función de `mapred.Process`, cuya firma no se ha incluido por simplicidad, controla el procesado de las instancias Hadoop a través de `Compute` y `Fabfile` (Fabric).

OpenStackBackend: es un pequeño backend que gestiona la conexión de los usuarios. Django incluye un sistema bastante completo para manejar las credenciales de conexión, pero en este proyecto hemos decidido apoyarnos en el de OpenStack para evitar la duplicidad de los usuarios, sus contraseñas y datos del perfil. Por tanto, el papel de este backend es hacer de *punte* entre el sistema integrado de credenciales de Django y el de OpenStack; cada petición de conexión se reenvía al *pipe* de acceso a OpenStack. Es decir, sólo es necesario que el usuario esté registrado en el cloud soporte, OpenStack en nuestro caso, para lanzar trabajos MapReduce.

OSServer y OSFlavor: son clases *helper* que desacoplan a la interfaz del servicio REST. Crean representaciones equivalentes en forma de objetos de los JSON que procesan como resultado de las invocaciones de `Compute`.

Diagrama de Clases — Objetos de Django

La figura 6.10 contiene, en detalle, las relaciones entre las *clases de objetos del modelo concreto* de Django.

django.db.models: este paquete lo proporciona la distribución de Django. Sirve de apoyo para escribir objetos del modelo de aplicación, basados en composiciones y extensiones de las clases contenidas.

Model: es la clase base de los objetos del modelo, esto es, todo objeto de nuestro minimundo ha de ser *descendiente* suyo. Usando esta clase y las definiciones de las clases de objetos, Django se encargará

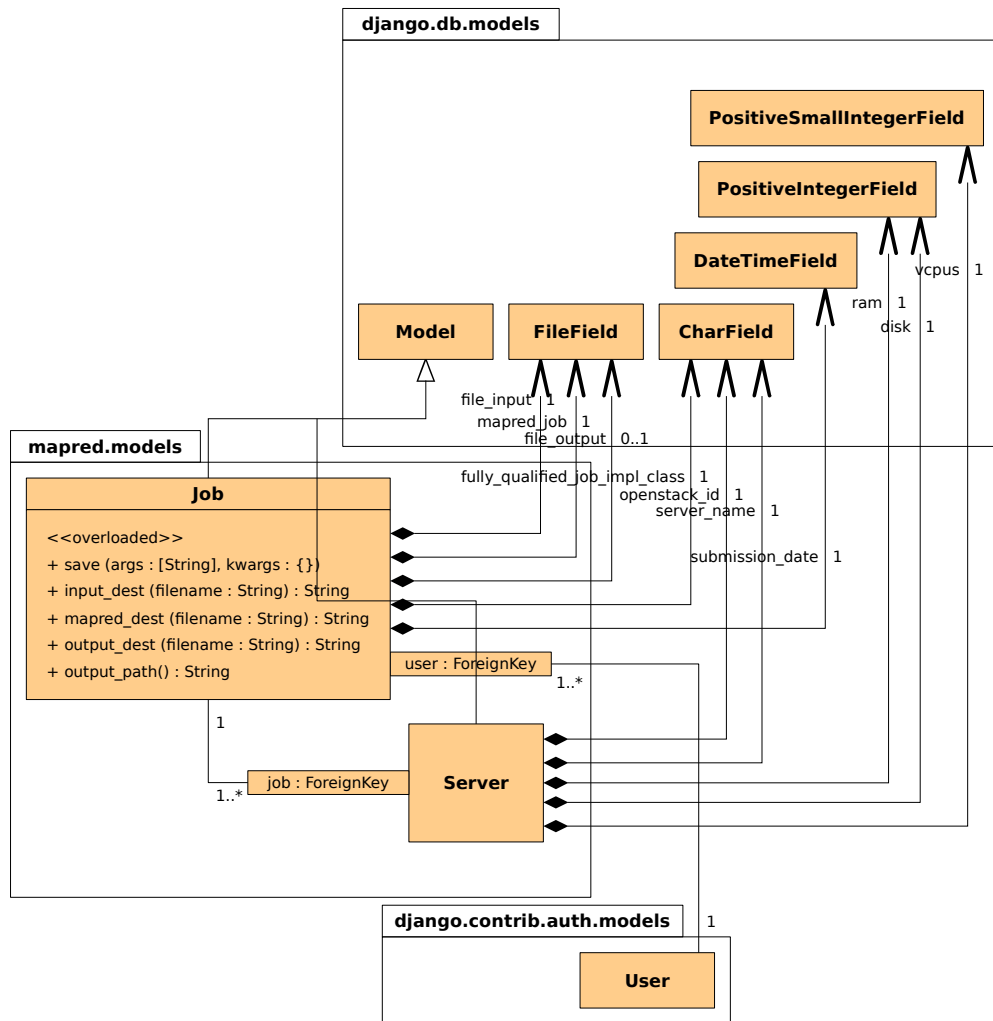


Figure 6.10: Diagrama de Clases — Objetos de Django

de gestionar todo lo relativo a la persistencia de los objetos del modelo en una base de datos.

Fields: son algunas de las implementaciones más usuales de campos de interés para los objetos que, junto con `Model`, permiten que Django establezca la morfología del esquema lógico de la base de datos asociada.

Job: representa un trabajo MapReduce para Hadoop. Comentar que se ha especializado la definición del método `save`, para poder organizar los trabajos en el sistema de ficheros por su identificador en la base de datos. Tal y como se ha comentado, la salvaguarda de información de cada `Job` en una base datos corre a cargo de Django.

Server: es la representación de la configuración individual de cada máquina virtual que participe en un trabajo Hadoop MapReduce.

User: es la clase interna de Django que almacena la información de los usuarios de la interfaz. Se utiliza como *Transfer Object* desde el sistema de autorización de OpenStack como portador de los datos de los usuarios en cada *sesión*.

Diagrama de Clases — Formularios de Django

La figura 6.11 muestra la descomposición en clases de los formularios que recogen las credenciales de acceso y la configuración del procesado. El paquete `django.forms` tiene organización y finalidad idénticos al paquete `django.db`.

`models` comentado anteriormente. En este caso `Form`, junto con los `Field` necesarios, es la clase extensible que aporta Django para concretar los formularios de usuario.

LoginForm: formulario de conexión de usuarios. Formado por un par de campos carácter que recogen el nombre de usuario y su contraseña en OpenStack.

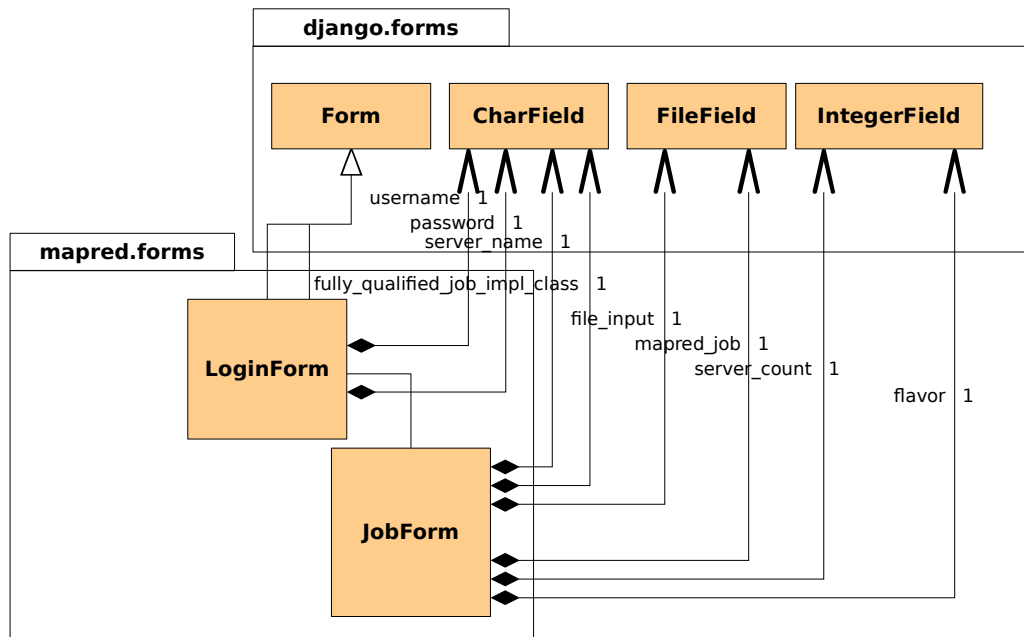


Figure 6.11: Diagrama de Clases — Formularios de Django

JobForm: es el formulario que permite definir la computación para MapReduce. Incluye: el prefijo del nombre de los servidores virtuales que se crearán, el nombre cualificado de la clase implementación de las funciones Map y Reduce, el fichero de entrada (paquete comprimido), el paquete *Jar* que contiene la clase implementación, el número de servidores virtuales necesarios y su *flavor* computacional.

Diagrama Entidad-Relación

Se había apuntado brevemente que Django posee la habilidad de construir automáticamente el esquema lógico de una base de datos relacional derivando el modelo de objetos escrito por el desarrollador. Era de esperar que la gestión de las operaciones *CRUD* (*Create, Read, Update, Delete*) sobre las tuplas almacenadas también estuviese implementada. La figura 6.12 muestra el diagrama Entidad-Relación correspondiente a la gestión de trabajos. Se han omitido las entidades y relaciones adyacentes que gestionan las sesiones

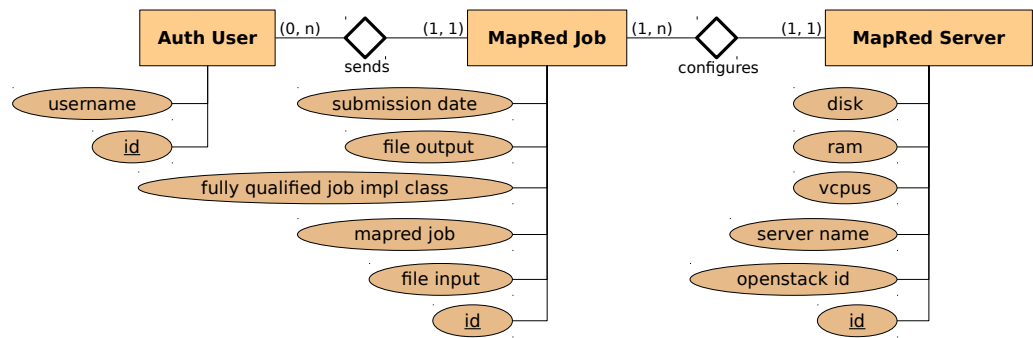


Figure 6.12: Diagrama Entidad-Relación

o las autorizaciones de los usuarios, entre otras, por transparencia.

Diagrama de Secuencia

Se presentan en las figuras 6.13 y 6.14 dos Diagramas de Secuencia. Reflejan el subconjunto más interesante de los mensajes intercambiados entre las entidades participantes en un procesado MapReduce. Estamos suponiendo que no se produce ningún error, que el usuario está conectado a la interfaz con sus credenciales, que introduce correctamente todos los datos de definición del procesado y que posee la autorización necesaria para lanzar máquinas virtuales en OpenStack. La figura 6.13 contiene la interacción completa. La 6.14 expone con detalle la activación posterior al mensaje #24 de la figura 6.13.

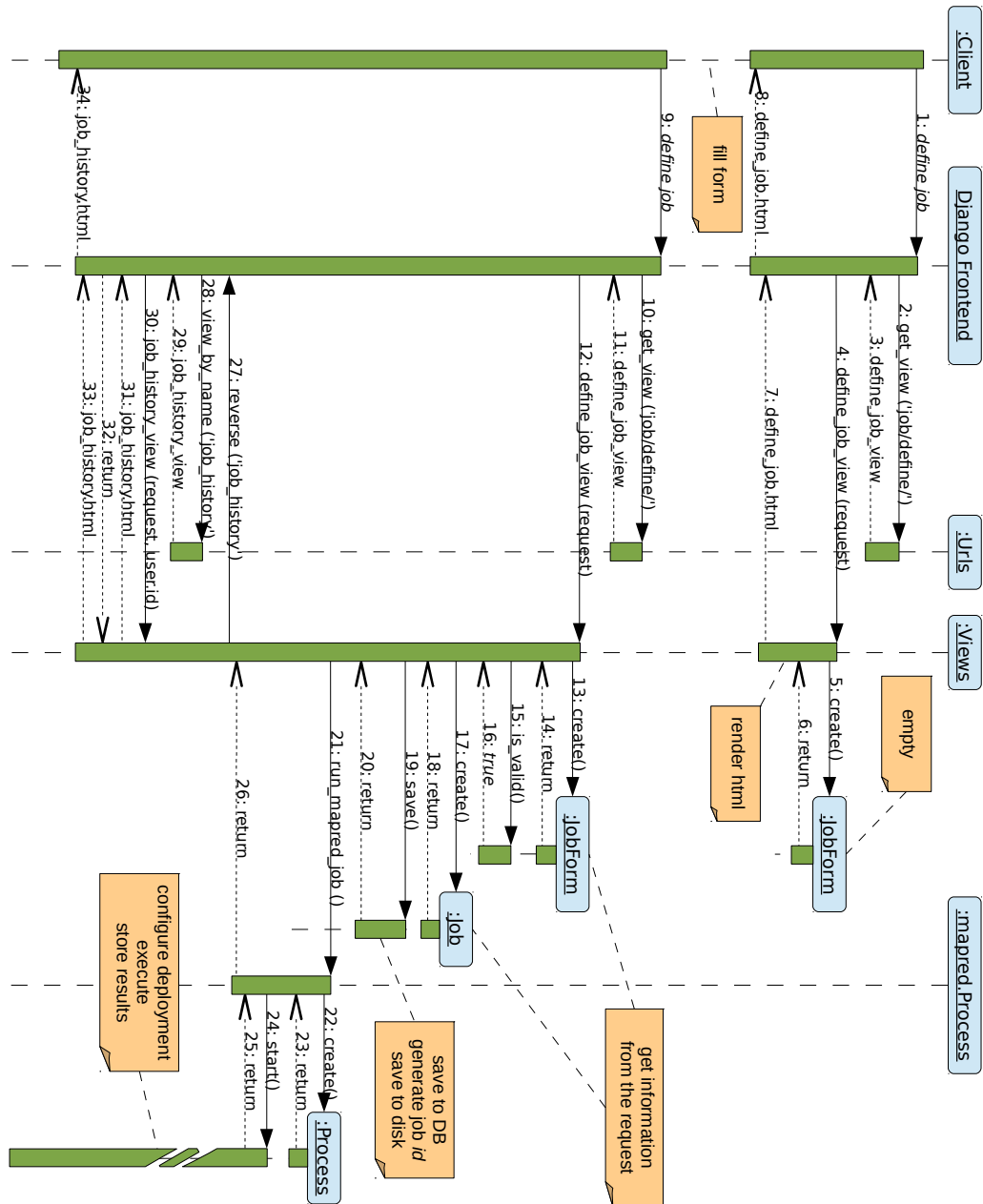


Figure 6.13: Diagrama de Secuencia (I)

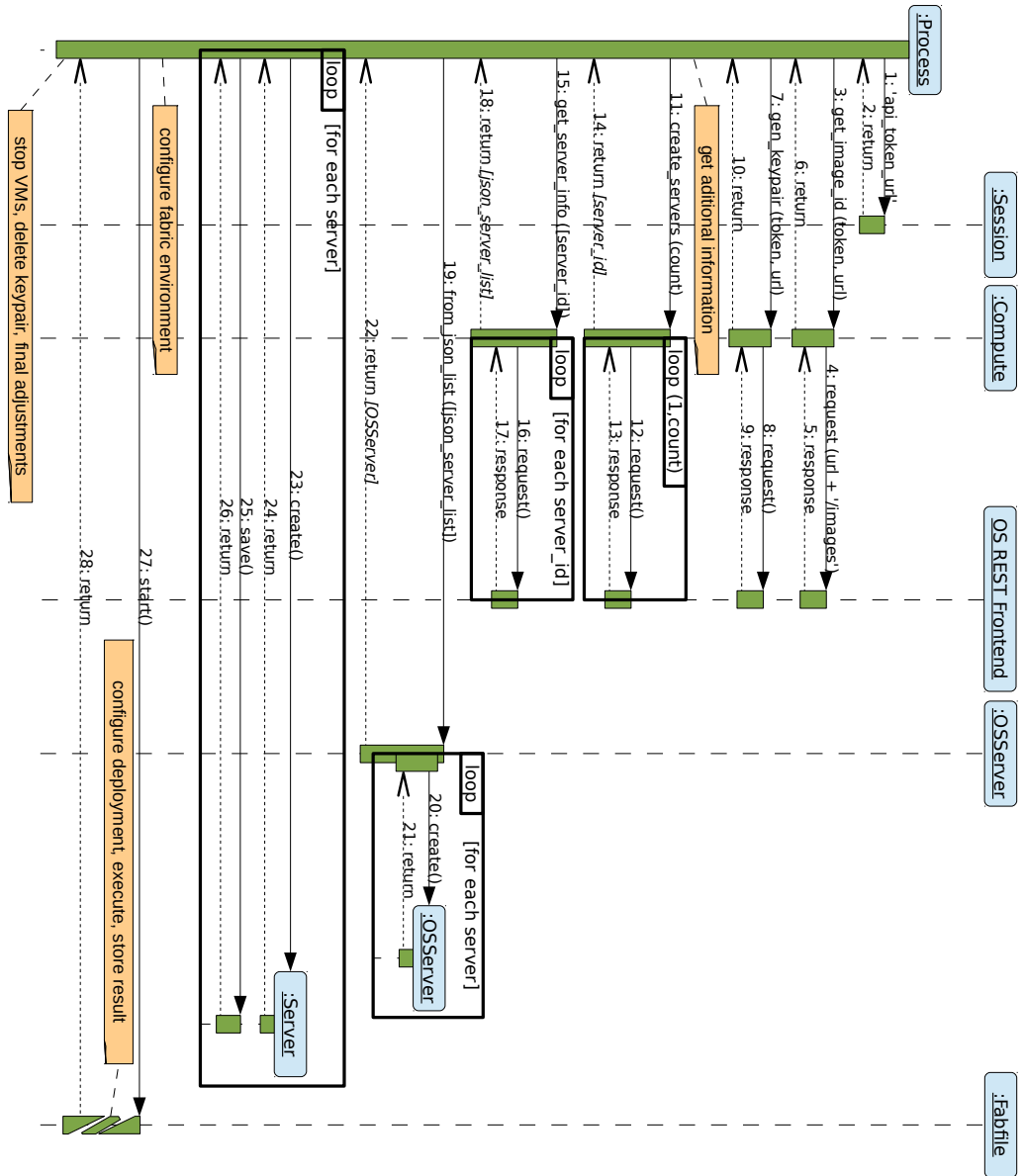


Figure 6.14: Diagrama de Secuencia (II)

Bibliography

- [1] Google Apps: Energy Efficiency in the Cloud. http://static.googleusercontent.com/external_content/untrusted_dlcp/www.google.com/en/us/green/pdf/google-apps.pdf. Accedido: junio 2013.
- [2] Jeffrey Dean and Sanjay Ghemawat. Mapreduce: Simplified Data Processing on Large Clusters. *Commun. ACM*, 51(1):107–113, 2008.
- [3] Pierre Riteau, Ancuta Iordache, and Christine Morin. Resilin: Elastic MapReduce for Private and Community Clouds. Research Report RR-7767, INRIA, October 2011.
- [4] Huan Liu and Dan Orban. Cloud MapReduce: A MapReduce Implementation on Top of a Cloud Operating System. In *Proceedings of the 2011 11th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing*, CCGRID '11, pages 464–474, Washington, DC, USA, 2011. IEEE Computer Society.
- [5] Steve Loughran, Jose Maria Alcaraz Calero, Andrew Farrell, Johannes Kirschnick, and Julio Guijarro. Dynamic Cloud Deployment of a MapReduce Architecture. *IEEE Internet Computing*, 16(6):40–50, November 2012.
- [6] Andy Edmonds, Thijs Metsch, Alexander Papaspyrou, and Alexis Richardson. Toward an Open Cloud Standard. *IEEE Internet Computing*, 16(4):15–25, July 2012.

-
- [7] Rich Uhlig, Gil Neiger, Dion Rodgers, Amy L. Santoni, Fernando C. M. Martins, Andrew V. Anderson, Steven M. Bennett, Alain Kagi, Felix H. Leung, and Larry Smith. Intel Virtualization Technology. *Computer*, 38(5):48–56, May 2005.
 - [8] Tom White. *Hadoop, the Definitive Guide*. O’ Reilly and Yahoo! Press, 2012.
 - [9] Nikita Ivanov. GridGain and Hadoop: Differences and Sinergies. <http://www.gridgain.com/blog/gridgain-hadoop-differences-synergies/>. Accedido: junio 2013.
 - [10] Jaliya Ekanayake, Hui Li, Bingjing Zhang, Thilina Gunarathne, Seung hee Bae, Judy Qiu, and Geoffrey Fox. Twister: A Runtime for Iterative MapReduce. In *The First International Workshop on MapReduce and its Applications*, 2010.
 - [11] Aaron McKenna, Matthew Hanna, Eric Banks, Andrey Sivachenko, Kristian Cibulskis, Andrew Kernytsky, Kiran Garimella, David Altshuler, Stacey Gabriel, Mark Daly, and Mark A. DePristo. The Genome Analysis Toolkit: A MapReduce Framework for Analyzing Next-generation DNA Sequencing Data. 2010.
 - [12] Quizmt project web page. <http://qizmt.myspace.com/>. Accedido: junio 2013.
 - [13] Adam Dou, Vana Kalogeraki, Dimitrios Gunopulos, Taneli Mielikainen, and Ville H. Tuulos. Misco: a Mapreduce Framework for Mobile Systems. In *Proceedings of the 3rd International Conference on Pervasive Technologies Related to Assistive Environments*, PETRA ’10, pages 32:1–32:8, New York, NY, USA, 2010. ACM.
 - [14] Peregrine project web page. http://peregrine_mapreduce.bitbucket.org/. Accedido: junio 2013.

-
- [15] Bingsheng He, Wenbin Fang, Qiong Luo, Naga K. Govindaraju, and Tuyong Wang. Mars: A MapReduce Framework on Graphics Processors. In *Proceedings of the 17th International Conference on Parallel Architectures and Compilation Techniques*, PACT '08, pages 260–269, New York, NY, USA, 2008. ACM.
 - [16] *Eucalyptus — 3.1.1 Installation Guide*, 2012.
 - [17] *CloudStack Basic Installation Guide for CloudStack Version 3.0.0 — 3.0.2*, 2012.
 - [18] Citrix Unveils Next Phase of Cloudstack Strategy. <http://www.citrix.com/news/announcements/apr-2012/citrix-unveils-next-phase-of-cloudstack-strategy.html>, 2012. Accedido: junio 2013.
 - [19] How to Use CloudStack without Hardware Virtualization. <http://support.citrix.com/article/CTX132015>, 2012. Accedido: junio 2013.
 - [20] *Apache CloudStack 4.0.0 — Incubating CloudStack Installation Guide*, 2012.
 - [21] *CloudStack Advanced Installation Guide for CloudStack Version 3.0.0 — 3.0.2*, 2012.
 - [22] *Citrix XenServer 6.0 Installation Guide*, 2012.
 - [23] OpenNebula 3.8.1 QuickStart. <http://wiki.centos.org/Cloud/OpenNebula/QuickStart>, 2012. Accedido: junio 2013.
 - [24] Getting Started with Openstack on Fedora 17. http://fedoraproject.org/wiki/Getting_started_with_OpenStack_on_Fedora_17, 2012. Accedido: junio 2013.
 - [25] *OpenStack Install and Deploy Manual — Red Hat — Folsom*, 2012.

-
- [26] *OpenStack Network (Quantum) Administration Guide — Folsom*, 2012.
 - [27] *OpenStack Object Storage Administration Manual — Folsom*, 2012.
 - [28] Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung. The Google File System. In *Proceedings of the Nineteenth ACM Symposium on Operating Systems Principles*, SOSP '03, pages 29–43, New York, NY, USA, 2003. ACM.
 - [29] *Amazon Elastic Compute Cloud — User Guide — API Version 2012-12-01*, 2013.
 - [30] QEMU Internals. <http://qemu.weilnetz.de/qemu-tech.html>, 2012. Accedido: junio 2013.
 - [31] Jaromír Hradílek, Douglas Silas, Martin Prpič, Stephen Wadeley, Eliška Slobodová, Tomáš Čapek, Petr Kovář, John Ha, David O'Brien, Michael Hideo, and Don Domingo. *Fedora 17 System Administrators' Guide. Deployment, Configuration and Administration of Fedora 17*. Red Hat Inc., 2012.
 - [32] Christopher Curran and Jan Mark Holzer. *Red Hat Enterprise Linux 5.2 – Virtualization Guide*. Red Hat Inc., 2008.
 - [33] Michael Hideo Smith. *Red Hat Enterprise Linux 5.2 — Deployment Guide*. Red Hat Inc., 2008.
 - [34] Johan De Gelas. Hardware Virtualization: the Nuts and Bolts. <http://www.anandtech.com/show/2480>, 2012. Accedido: junio 2013.
 - [35] *OpenStack Install and Deploy Manual — Red Hat — Essex*, 2012.
 - [36] *OpenStack Compute Administration Manual — Essex*, 2012.
 - [37] *OpenStack Compute Administration Manual — Folsom*, 2012.
 - [38] Jacek Artymiak. *Programming OpenStack Compute API*. Rackspace US Inc., 2012.

-
- [39] *OpenStack API Reference*, 2013.
 - [40] OpenNebula OCCI Specification 3.8. <http://opennebula.org/documentation:archives:rel3.8:occidd>, 2012. Accedido: junio 2013.
 - [41] DEISA. Deisa glossary. <http://www.deisa.eu/references>. Accedido: junio 2013.
 - [42] Peter Mell and Timothy Grance. The NIST Definition of Cloud Computing. <http://csrc.nist.gov/publications/nistpubs/800-145/SP800-145.pdf>, 2011. Accedido: junio 2013.
 - [43] *CloudStack API Documentation (v3.0)*, 2012.
 - [44] Simon Kelley. dnsmasq — A Lightweight DHCP and Caching DNS Server — ManPage. <http://www.thekelleys.org.uk/dnsmasq/docs/dnsmasq-man.html>, 2012. Accedido: junio 2013.
 - [45] HDFS Users' Guide. http://hadoop.apache.org/docs/r1.0.4/single_node_setup.html, 2010. Accedido: junio 2013.
 - [46] Single Node Setup. http://hadoop.apache.org/docs/r1.0.4/single_node_setup.html, 2008. Accedido: junio 2013.
 - [47] MapReduce Tutorial. http://hadoop.apache.org/docs/r1.0.4/mapred_tutorial.html, 2008. Accedido: junio 2013.
 - [48] Cluster setup. http://hadoop.apache.org/docs/r1.0.4/hdfs_user_guide.html, 2008. Accedido: junio 2013.
 - [49] Dhruba Borthakur. HDFS Architecture Guide. http://hadoop.apache.org/docs/r1.0.4/hdfs_design.html, 2012. Accedido: junio 2013.
 - [50] Jimmy Lin and Chris Dyer. *Data-Intensive Text Processing with MapReduce*. 2010.

- [51] Sriram Rao, Raghu Ramakrishnan, Adam Silberstein, Mike Ovsianikov, and Damian Reeves. Sailfish: A Framework for Large Scale Data Processing. In *Proceedings of the Third ACM Symposium on Cloud Computing*, SoCC '12, pages 4:1–4:14, New York, NY, USA, 2012. ACM.
- [52] Ahmed Metwally and Christos Faloutsos. V-SMART-join: A Scalable MapReduce Framework for All-pair Similarity Joins of Multisets and Vectors. *Proc. VLDB Endow.*, 5(8):704–715, April 2012.
- [53] Amr Awadallah. Apache Hadoop in the Enterprise — Keynote. Cloudera Inc., 2011.
- [54] Mendel Cooper. Advanced Bash-Scripting Guide. An In-depth Exploration of the Art of Shell Scripting. <http://tldp.org/LDP/abs/html/>, 2012. Accedido: junio 2013.
- [55] Bruce Barnett. Sed — An Introduction and Tutorial. <http://www.grymoire.com/Unix/Sed.html>, 2012. Accedido: junio 2013.
- [56] MapReduce Wikipedia entry. <http://en.wikipedia.org/wiki/MapReduce>, 2012. Accedido: junio 2013.