
Final Year Internship Project

*Presented to obtain a
Computer Science Engineer Diploma*

Specialization : Information System

Subject :

**Development and Design of a Community and social
platform for Orange**

Made by

KAIS NEFFATI

Host company : SOFRECOM TUNISIE



Entreprise supervisors : Mr Walid HAMOUDA AND Mr Mehdi KALAI

University supervisor : Mrs Ilhem Abdelmouleh

JUIN 2017

College year : 2016/2017

ABSTRACT

Kais NEFFATI, Computer Science Engineering , National Engineering School of Carthage Tunisia

Abstract of a computer science engineer final year internship project , Submitted 29 Juin 2017

Development and Design of a Community and social platform for Orange

The aim of this project is to design and develop a customized marketing system that allow Data analytics. Today, This kind of systems are widely used. And have gained so much followers. These systems have successfully been one of the most powerful marketing strategies. Today, the Orange group as a world wide company wants to implement one of these solutions. For this reason the author (me) decided to commit to this project and consider it as its graduation project. This thesis first presents the host company, put the project in context and introduce the development, analysis and design methodologies. In a second stage, it introduces several technical and functional requirements specifications. Then it exposes the system architectural and functional design. Finally, we implement and expose our system features.

DEDICATION AND ACKNOWLEDGEMENTS

Throughout my life my family has always been there during those difficult and trying times. I would like to dedicate this thesis and everything I do to my father Noureddine NEFFATI, my mother Yosra REZGUI, my sister Chaima NEFFATI, my little sister Nesrine NEFFATI my brother-in-law Helmi CHEBBI and especially my nephew Zakaria CHEBBI. Thanks to all my friends for sharing my happiness when starting this project and following with encouragement when it seemed too difficult to be completed. I would have probably given up without their support and example on what to do when you really want something. At the end of this work, I would like to express my sincere thanks to all the people Which have contributed, either directly or indirectly, to its smooth operation. I would like to express my gratitude and sincere thanks to Mrs Ilhem ABDELHEDI, My supervisor at Enicarthage, who knew with his professionalism and his pedagogical supervision Guide me during these four months of training. I thank her for her regular follow-up My advancement in the project, for his valuable advice and encouragement. I would also like to thank my supervisors at Sofrecom Tunisia, Mr Walid HAMOUDA, Development engineer, and Mr Mehdi KALAI, technical manager, for their availability, Their advice and confidence. I would also like to thank the entire Sofrecom Tunisia team for their hospitality, support and Their team spirit.

AUTHOR'S DECLARATION

I declare that the work in this dissertation was carried out in accordance with the requirements of the University's Regulations and Code of Practice for Project Degree Programmes and that it has not been submitted for any other academic award. Except where indicated by specific reference in the text, the work is the candidate's own work. Work done in collaboration with, or with the assistance of, others, is indicated as such. Any views expressed in the dissertation are those of the author.

KAIS NEFFATI

SIGNED: DATE:

ILHEM ADBELHEDI

SIGNED: DATE:

WALID HAMOUDA

SIGNED: DATE:

MEHDI KALAI

SIGNED: DATE:

TABLE OF CONTENTS

	Page
List of Tables	ix
List of Figures	xi
1 Project Scope	1
1.1 Host company	1
1.1.1 Sofrecom	1
1.1.2 Sofrecom Tunisie	2
1.2 Problematic	3
1.3 Solution	3
1.3.1 Intorduction	3
1.3.2 Big Data Architecture	3
1.3.3 Application Architecture	5
1.3.4 Synthesis	10
1.4 Development methodology	10
1.4.1 Introduction	10
1.4.2 Implementation	11
1.5 Analysis and design methodology	13
1.5.1 Introduction	13
1.5.2 Abstraction	13
1.5.3 Models VS. Methodologies	14
1.5.4 UML 2.0	15
1.6 Conclusion	15
2 Requirements Analysis and Specification	17
2.1 Functional branch	18
2.1.1 Requirements analysis	18
2.1.2 Requirements specification	19
2.2 Technical branch	26
2.2.1 Big Data architecture	26

TABLE OF CONTENTS

2.2.2	Architecture microservice	39
2.2.3	Deployment architecture	43
2.3	Conclusion	44
3	Design	45
3.1	General design	45
3.1.1	Architectural design	45
3.2	Detailed design	49
3.2.1	Package Diagram	49
3.2.2	Class Diagram	49
3.2.3	Object Sequence Diagrams	50
3.3	Conclusion	53
4	Implementation	55
4.1	Work environnement	55
4.1.1	Hardware environnement	55
4.2	Software environnement	55
4.3	Realized work	56
4.3.1	Environnement preparation	56
4.3.2	Log Producer	61
4.3.3	Spark Jobs	62
4.3.4	Show Resualts on Apache Zeppelin	68
4.3.5	Account Managment	68
4.3.6	Promotion Managment	70
4.3.7	Display Analytics	71
4.3.8	Heat Map users	72
4.3.9	Summary	72
4.4	General conclusion	73
A	Appendix A	75
Bibliography		77

LIST OF TABLES

TABLE	Page
--------------	-------------

LIST OF FIGURES

FIGURE	Page
1.1 filiales et clients de Sofrecom.	1
1.2 filiales et clients de Sofrecom.	2
1.3 filiales et clients de Sofrecom.	4
1.4 filiales et clients de Sofrecom.	5
1.5 filiales et clients de Sofrecom.	6
1.6 filiales et clients de Sofrecom.	9
1.7 filiales et clients de Sofrecom.	11
1.8 filiales et clients de Sofrecom.	12
2.1 filiales et clients de Sofrecom.	17
2.2 filiales et clients de Sofrecom.	20
2.3 filiales et clients de Sofrecom.	21
2.4 filiales et clients de Sofrecom.	22
2.5 filiales et clients de Sofrecom.	24
2.6 filiales et clients de Sofrecom.	25
2.7 filiales et clients de Sofrecom.	26
2.8 filiales et clients de Sofrecom.	27
2.9 filiales et clients de Sofrecom.	27
2.10 filiales et clients de Sofrecom.	28
2.11 filiales et clients de Sofrecom.	28
2.12 filiales et clients de Sofrecom.	29
2.13 filiales et clients de Sofrecom.	29
2.14 filiales et clients de Sofrecom.	30
2.15 filiales et clients de Sofrecom.	31
2.16 filiales et clients de Sofrecom.	31
2.17 filiales et clients de Sofrecom.	35
2.18 filiales et clients de Sofrecom.	35
2.19 filiales et clients de Sofrecom.	36
2.20 filiales et clients de Sofrecom.	38

LIST OF FIGURES

2.21 filiales et clients de Sofrecom	41
3.1 filiales et clients de Sofrecom	46
3.2 filiales et clients de Sofrecom	46
3.3 filiales et clients de Sofrecom	47
3.4 filiales et clients de Sofrecom	47
3.5 filiales et clients de Sofrecom	48
3.6 filiales et clients de Sofrecom	49
3.7 filiales et clients de Sofrecom	50
3.8 filiales et clients de Sofrecom	51
3.9 filiales et clients de Sofrecom	52
4.1 filiales et clients de Sofrecom	57
4.2 filiales et clients de Sofrecom	57
4.3 filiales et clients de Sofrecom	58
4.4 filiales et clients de Sofrecom	58
4.5 filiales et clients de Sofrecom	59
4.6 filiales et clients de Sofrecom	59
4.7 filiales et clients de Sofrecom	60
4.8 filiales et clients de Sofrecom	61
4.9 filiales et clients de Sofrecom	61
4.10 filiales et clients de Sofrecom	62
4.11 filiales et clients de Sofrecom	63
4.12 filiales et clients de Sofrecom	64
4.13 filiales et clients de Sofrecom	64
4.14 filiales et clients de Sofrecom	65
4.15 filiales et clients de Sofrecom	65
4.16 filiales et clients de Sofrecom	66
4.17 filiales et clients de Sofrecom	66
4.18 filiales et clients de Sofrecom	67
4.19 filiales et clients de Sofrecom	67
4.20 filiales et clients de Sofrecom	68
4.21 filiales et clients de Sofrecom	68
4.22 filiales et clients de Sofrecom	69
4.23 filiales et clients de Sofrecom	69
4.24 filiales et clients de Sofrecom	70
4.25 filiales et clients de Sofrecom	70
4.26 filiales et clients de Sofrecom	71
4.27 filiales et clients de Sofrecom	71

LIST OF FIGURES

4.28 filiales et clients de Sofrecom.	72
---	----

PROJECT SCOPE

In this chapter, we introduce the host company "Sofrecom Tunisie". Then we explain the project by putting it in context , exposing the problematic and proposing a solution. Finally, we describe the development methodologies and the analysis/Design one.

1.1 Host company

1.1.1 Sofrecom

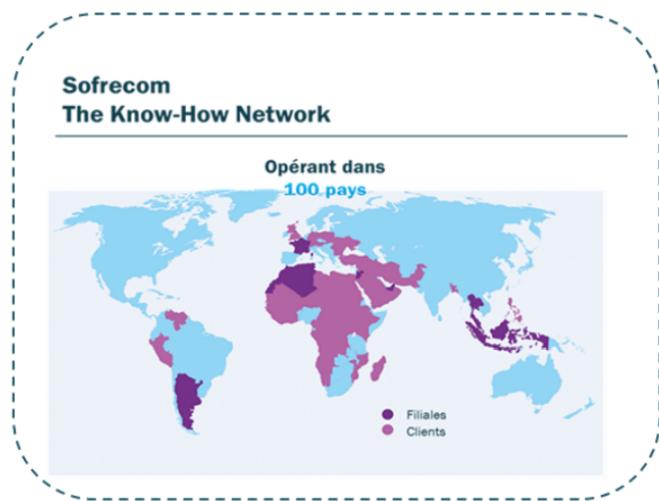


FIGURE 1.1. Sofrecom's subsidiaries and customers. Figure taken from [1].

Sofrecom, an Orange subsidiary, has developed over 50 years unique know-how about oper-

ators businesses, making it a world leading specialist in telecommunications consultancy and engineering. Its experience of mature and emerging markets, combined with its deep understanding of the structuring changes affecting the telecoms market, make it a valued partner for operators, governments and international investors. In recent years, 200 major players in over 100 countries have entrusted strategic and operational projects to Sofrecom. Sofrecom is above all a network of men and women, a powerful network of know-how and expertise that ties its personnel to customers, Orange experts and industrial and local partners. Sofrecom's Know-How Network is also the guarantee of effective transfer of knowhow and skills for sustainable transformation based on internationally certified methodologies.

1.1.2 Sofrecom Tunisie



FIGURE 1.2. Sofrecom's subsidiaries and customers. Figure taken from [1].

Sofrecom Tunisie, a young company in the local market, inaugurated in October 2012, is the youngest and largest subsidiary of the Sofrecom group in the Africa and Middle East zone. In 5 years it has positioned itself as a major player in consulting and engineering in telecommunications. More than 400 Experts, Sofrecom Tunisia has acquired its strength and credibility through hundreds of diversified international projects allowing us to be one step ahead of our competitors. As part of the Orange Group, with an international reputation, working for two major Orange customers: DSI France and OLS

1.2 Problematic

When we talk about selling products and increasing revenues, we totally refer to making offers and publishing promotions. So many companies are using a variety of ways to promote their service. One of them is the Orange group. Orange Group is spending a lot of money on marketing campaigns. Orange is using SMS, MMS and TV emission to be attached to its clients which is not always the best way. As many marketing experts said, the key success of a marketing strategy is to know where, when and to whom a product and/or a service should be exposed.

A team of engineers in Sofrecom Tunisia, and in an innovation scope, proposes to make a new dimension in Orange marketing strategy. The idea is to provide a mobile platform that uses proximity to allow users to see promotions, services, Events, Sensitization ... depending on their interests and their locations. This application will have a free, location-based orange iBeacons, use device's GPS capability to locate Orange iBeacons. Subscribers can capture gifts and promotions when nearby iBeacon using BLE devices through application. The Backend engine of this application gets customers relevant information (rate plan, balance, services,...) in order to propose a bespoke award. This idea brings some added values; promote Orange low-rated services usage, increase customers loyalty, trigger based promotions and gifts as per user preferences and targeting subscribers based on usages and consumption habits (Voices, SMS, Data, TopUp...).

Bringing this application to life requires implementing a system that allows us to recommend a service and track user behavior. This system should be able to deal with a huge amount of data (Orange Clients), a Big data should be analyzed. We were influenced the most by YouTube, which, in first use, exposes videos in an arbitrary way. As users begin to interact with the dashboard YouTube recommends the right list of videos. YouTube takes into consideration recent and old user data. Once the user changes his behavior YouTube will try to rebalance the dashboard view. Our system finally should be able to analyze big data, make recommendations based on interests and proximity and support the logic of application in the backend.

1.3 Solution

1.3.1 Introduction

Our application is divided mainly in two parts, a part responsible for data analytics and another one for application logic. Many patterns and types of architecture are required.

1.3.2 Big Data Architecture

1.3.2.1 Lambda Architecture

is a data-processing architecture designed to handle massive quantities of data by taking advantage of both batch- and stream-processing methods. This approach to architecture attempts to

balance latency, throughput, and fault-tolerance by using batch processing to provide comprehensive and accurate views of batch data, while simultaneously using real-time stream processing to provide views of online data. The two view outputs may be joined before presentation. The rise of lambda architecture is correlated with the growth of big data, real-time analytics, and the drive to mitigate the latencies of map-reduce. Lambda architecture depends on a data model with an append-only, immutable data source that serves as a system of record. It is intended for ingesting and processing timestamped events that are appended to existing events rather than overwriting them. State is determined from the natural time-based ordering of the data.

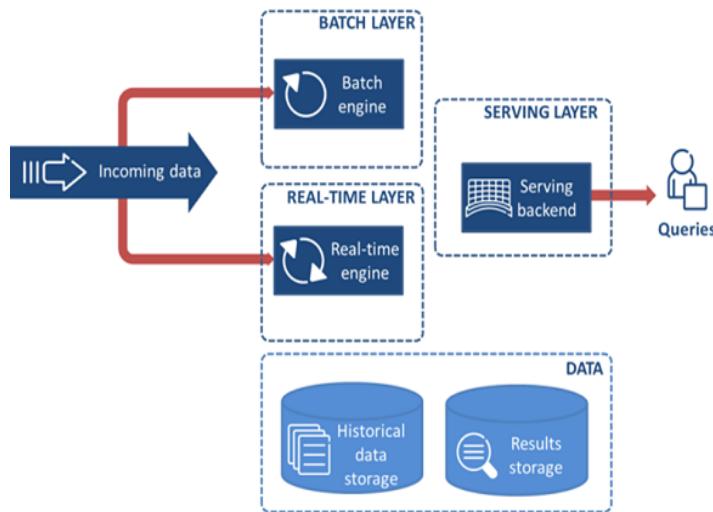


FIGURE 1.3. Sofrecom's subsidiaries and customers. Figure taken from [1].

Batch Layer Precomputes results using a distributed processing system that can handle very large quantities of data. The batch layer aims at perfect accuracy by being able to process all available data when generating views. This means it can fix any errors by recomputing based on the complete data set, then updating existing views. Output is typically stored in a read-only database, with updates completely replacing existing precomputed views.

Streaming Layer The speed layer processes data streams in real time and without the requirements of fix-ups or completeness. This layer sacrifices throughput as it aims to minimize latency by providing real-time views into the most recent data. Essentially, the speed layer is responsible for filling the "gap" caused by the batch layer's lag in providing views based on the most recent data. This layer's views may not be as accurate or complete as the ones eventually produced by the batch layer, but they are available almost immediately after data is received, and can be replaced when the batch layer's views for the same data become available.

Serving Layer Output from the batch and speed layers are stored in the serving layer, which responds to ad-hoc queries by returning precomputed views or building views from the processed data.

1.3.2.2 Kappa Architecture

is a simplification of Lambda Architecture. A Kappa Architecture system is like a Lambda Architecture system with the batch processing system removed. To replace batch processing, data is simply fed through the streaming system quickly.

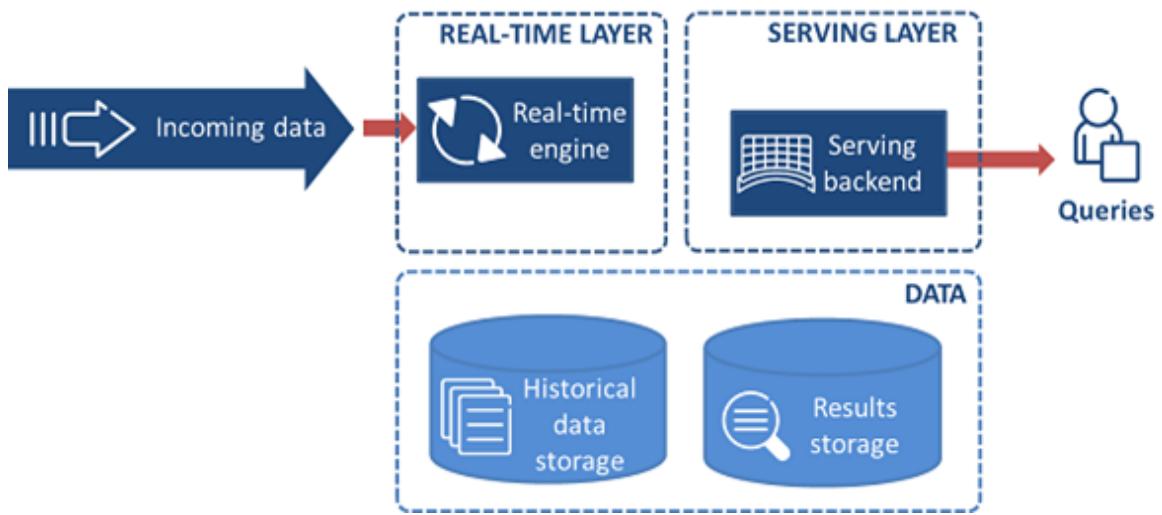


FIGURE 1.4. Sofrecom's subsidiaries and customers. Figure taken from [1].

1.3.2.3 Synthesis

The question now is which one fit the best for us? In fact, choosing whether to work with lambda or kappa architecture depends on the use case, so if we treat the same way streaming and existing data the kappa architecture would be a great solution. And if we treat them differently, we should implement the lambda one. In our case, promotions will be recommended upon historical and streaming data. So the best architecture for our case is the lambda architecture. Still, one thing, Which tools shall we use!

1.3.3 Application Architecture

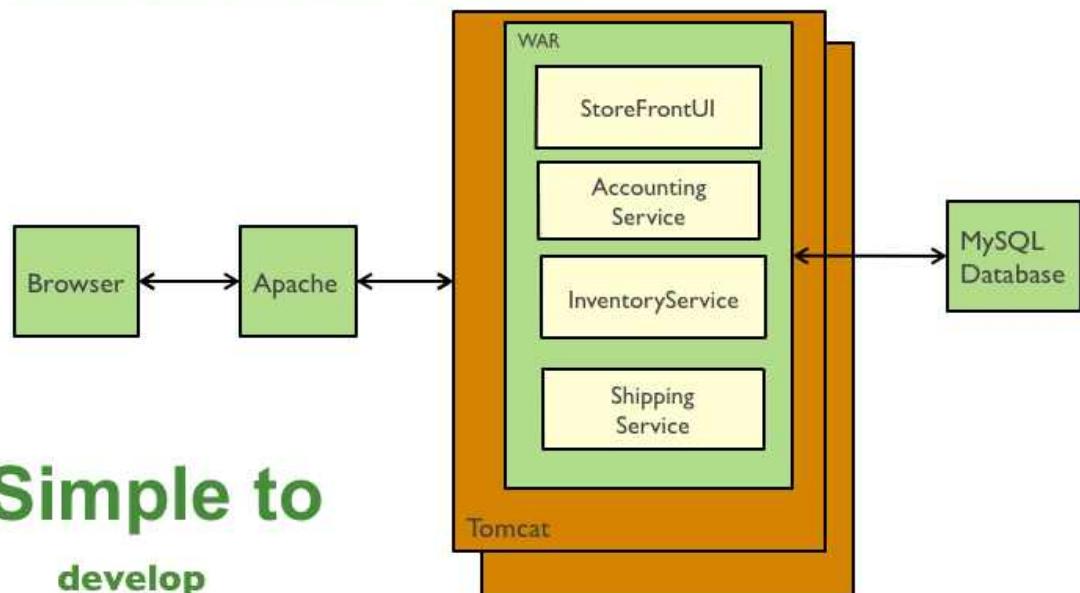
we are developing a server-side enterprise application that handle the logic of application. It must support a variety of different clients including desktop browsers, mobile browsers and native mobile applications. This application might also expose some APIs for 3rd parties to consume. It

might also integrate with other applications via either web services or a message broker. The application handles requests (HTTP requests and messages) by executing business logic; accessing a database; exchanging messages with other systems; and returning a HTML/JSON/XML response. There are logical components corresponding to different functional areas of the application.

1.3.3.1 Monolithic application

One of the solutions is to build an application with a monolithic architecture. The application consists of several components and deployed as a single monolithic application. For example, a Java web application consists of a single WAR file that runs on a web container such as Tomcat. A Rails application consists of a single directory hierarchy deployed using either, for example, Phusion Passenger on Apache/Nginx or JRuby on Tomcat. we can run multiple instances of the application behind a load balancer in order to scale and improve availability.

Traditional web application architecture



Simple to
develop
test
deploy
scale

FIGURE 1.5. Sofrecom's subsidiaries and customers. Figure taken from [1].

This solution has a number of benefits:

- Simple to develop - the goal of current development tools and IDEs is to support the development of monolithic applications

- Simple to deploy - you simply need to deploy the WAR file (or directory hierarchy) on the appropriate runtime
- Simple to scale - you can scale the application by running multiple copies of the application behind a load balancer

However, once the application becomes large and the team grows in size, this approach has a number of drawbacks that become increasingly significant:

- The large monolithic code base intimidates developers, especially ones who are new to the team. The application can be difficult to understand and modify. As a result, development typically slows down. Also, because there are not hard module boundaries, modularity breaks down over time. Moreover, because it can be difficult to understand how to correctly implement a change the quality of the code declines over time. It's a downwards spiral.
- Overloaded IDE - the larger the code base the slower the IDE and the less productive developers are.
- Overloaded web container - the larger the application the longer it takes to start up. This had have a huge impact on developer productivity because of time wasted waiting for the container to start. It also impacts deployment too.
- Continuous deployment is difficult - a large monolithic application is also an obstacle to frequent deployments. In order to update one component you have to redeploy the entire application. This will interrupt background tasks (e.g. Quartz jobs in a Java application), regardless of whether they are impacted by the change, and possibly cause problems. There is also the chance that components that haven't been updated will fail to start correctly. As a result, the risk associated with redeployment increases, which discourages frequent updates. This is especially a problem for user interface developers, since they usually need to iterate rapidly and redeploy frequently.
- Scaling the application can be difficult - a monolithic architecture is that it can only scale in one dimension. On the one hand, it can scale with an increasing transaction volume by running more copies of the application. Some clouds can even adjust the number of instances dynamically based on load. But on the other hand, this architecture can't scale with an increasing data volume. Each copy of application instance will access all of the data, which makes caching less effective and increases memory consumption and I/O traffic. Also, different application components have different resource requirements - one might be CPU intensive while another might memory intensive. With a monolithic architecture we cannot scale each component independently

- Obstacle to scaling development - A monolithic application is also an obstacle to scaling development. Once the application gets to a certain size it's useful to divide up the engineering organization into teams that focus on specific functional areas. For example, we might want to have the UI team, accounting team, inventory team, etc. The trouble with a monolithic application is that it prevents the teams from working independently. The teams must coordinate their development efforts and redeployments. It is much more difficult for a team to make a change and update production.
- Requires a long-term commitment to a technology stack - a monolithic architecture forces you to be married to the technology stack (and in some cases, to a particular version of that technology) you chose at the start of development. With a monolithic application, can be difficult to incrementally adopt a newer technology. For example, let's imagine that you chose the JVM. You have some language choices since as well as Java you can use other JVM languages that inter-operate nicely with Java such as Groovy and Scala. But components written in non-JVM languages do not have a place within your monolithic architecture. Also, if your application uses a platform framework that subsequently becomes obsolete then it can be challenging to incrementally migrate the application to a newer and better framework. It's possible that in order to adopt a newer platform framework you have to rewrite the entire application, which is a risky undertaking.

1.3.3.2 Microservices Architecture

another solution is to define an architecture that structures the application as a set of loosely coupled, collaborating services. This approach corresponds to the Y-axis of the Scale Cube. Each service implements a set of narrowly, related functions. For example, an application might consist of services such as the order management service, the customer management service etc. Services communicate using either synchronous protocols such as HTTP/REST or asynchronous protocols such as AMQP. Services can be developed and deployed independently of one another. Each service has its own database in order to be decoupled from other services. Data consistency between services is maintained using an event-driven architecture .

This solution has a number of benefits:

- Each microservice is relatively small (Easier for a developer to understand, The IDE is faster making developers more productive, The application starts faster, which makes developers more productive, and speeds up deployments)
- Each service can be deployed independently of other services - easier to deploy new versions of services frequently
- Easier to scale development. It enables you to organize the development effort around multiple teams. Each (two pizza) team is owns and is responsible for one or more single

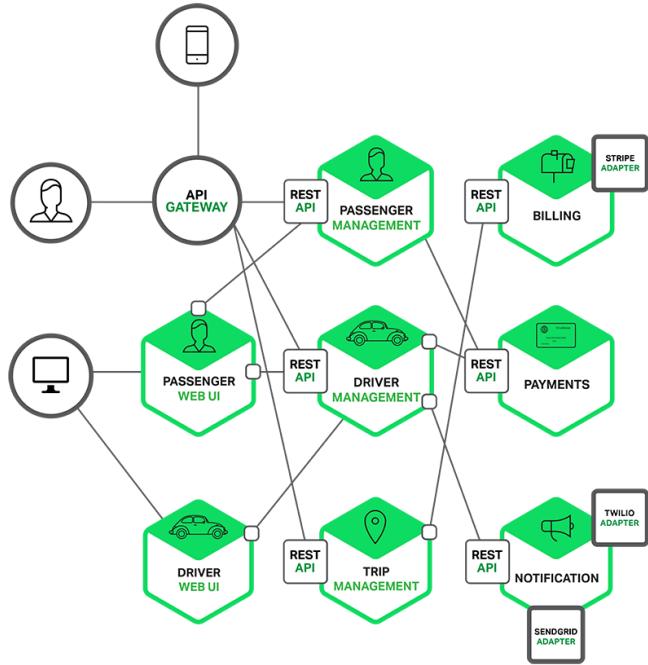


FIGURE 1.6. Sofrecom's subsidiaries and customers. Figure taken from [1].

service. Each team can develop, deploy and scale their services independently of all of the other teams.

- Improved fault isolation. For example, if there is a memory leak in one service then only that service will be affected. The other services will continue to handle requests. In comparison, one misbehaving component of a monolithic architecture can bring down the entire system.
- Each service can be developed and deployed independently
- Eliminates any long-term commitment to a technology stack. When developing a new service you can pick a new technology stack. Similarly, when making major changes to an existing service you can rewrite it using a new technology stack.

This solution has a number of drawbacks:

- Developers must deal with the additional complexity of creating a distributed system (Developer tools/IDEs are oriented on building monolithic applications and don't provide explicit support for developing distributed applications. Testing is more difficult , Developers must implement the inter-service communication mechanism , Implementing use cases that span multiple services without using distributed transactions is difficult , Implementing use cases that span multiple services requires careful coordination between the teams)

- Deployment complexity. In production, there is also the operational complexity of deploying and managing a system comprised of many different service types.
- Increased memory consumption. The microservice architecture replaces N monolithic application instances with NxM services instances. If each service runs in its own JVM (or equivalent), which is usually necessary to isolate the instances, then there is the overhead of M times as many JVM runtimes. Moreover, if each service runs on its own VM (e.g. EC2 instance), as is the case at Netflix, the overhead is even higher.

1.3.4 Synthesis

Choosing whether to work with microservices or monolithic application is depending on the project maturity. So when developing the first version of an application, we often do not have the problems that the microservice approach solves. Moreover, using an elaborate, distributed architecture will slow down development. But looking to the benefits that offer the microservices architecture is really impressive. And as this project is dedicated to a big company with a big customer range. This application should be apt to scale up quickly and ready to add so much features. So implementing the microservices architecture is a great choice. The question now is which tools or patterns shall we use ?

1.4 Development methodology

1.4.1 Introduction

2TUP is a unified process (i.e. a software development process) built on the UML modeling language. Every process answers the following main characteristics:

- It is an incremental process, allowing a better technical and functional risk management and thus constituting the deadlines and the costs control.
- It is an iterative process. The degrees of abstraction are increasingly precise at each iteration.
- It is component oriented, offering flexibility to the model and supporting the re-use.
- It is user oriented because built from their expectations.

The 2TUP process answers to the constraints of change of the information systems subjected themselves to two types of constraints: functional constraints and technical constraints as shows it the following diagram: In concrete terms, the process is modeled by two branches (tracks):

- A functional track (capitalization of knowledge trade)
- A technical track (re-use of a technical knowhow).



FIGURE 1.7. Sofrecom's subsidiaries and customers. Figure taken from [1].

Then these two tracks amalgamate for the realization of the system. This is why this process is still called Yshaped process. With this development process, a model is essential in order to anticipate the results. A model can be used with each step of the development with an increasing detailed manner. The industrial standard of object modeling, UML, was selected as the development tool. It appeared very difficult to consider the 2TUP process without using UML and, more particularly UML 2.0 which support the oriented design component.

1.4.2 Implementation

1.4.2.1 Introduction

The first phase of the implementation starts with the preliminary study which introduces the project. The specifications, initial document of the functional and technical needs, are partly the result of this work. It is supplemented by the modeling of the total system context. Not everything is to be described at this stage but simply to identify the external entities interacting (actors), to list the interactions (messages) and to represent this unit on a model (context).

1.4.2.2 Feature track

The functional branch makes an inventory of the functional needs and analyzes it. This phase formalizes and specifies the elements of the preliminary study. The applied use case technique translates the whole interactions between the system and the actors. The obtained use cases are then organized (treated on a hierarchical basis, generalized, specialized...). They make it possible to identify the classes and they permit the oriented object modeling generated in the analysis part.

1.4.2.3 Technical track

The technical branch lists the technical needs and proposes a generic design validated by a prototype. The pre-necessary techniques revealed in the preliminary study, showing the operational needs and the strategic choices of development, lead to the development of the construction process. To do this, several stages are necessary:

- The inventory of technical specifications related to the hardware.

- The inventory of the software specifications.

1.4.2.4 Middle track

The medium branch supports the preliminary design, the detailed design, coding, the tests and the validation. The preliminary design is one of the most sensitive steps of the 2TUP process. It represents the fusion of the functional and technical tracks. It finishes when the deployment model (working stations, architectures), the operating model (components, applications), the logical model (classes diagrams, classes), interfaces (users and components) and the software configuration model are defined.

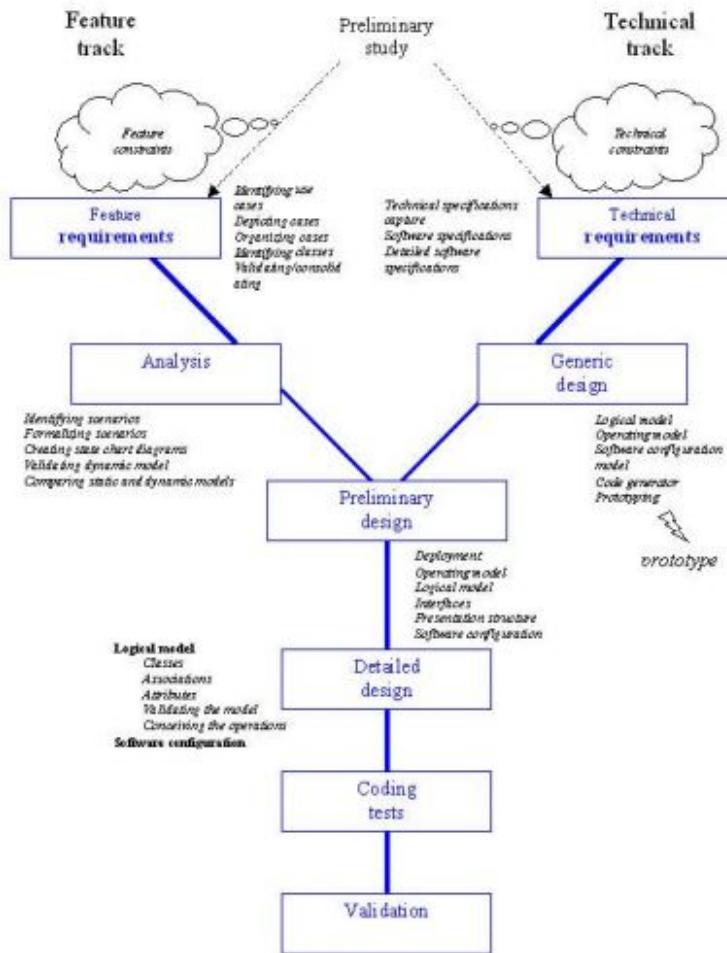


FIGURE 1.8. Sofrecom's subsidiaries and customers. Figure taken from [1].

The detailed design conceives and documents very exactly the code that will be generated. It is largely founded on UML representations and implements, in an iterative manner, a system construction process to obtain a "model ready to code".

The various components carried out in the detailed design are coded. The code units obtained are tested as they are written and produced. The validation, called ,*réponse*, step, consists of the homologation of the developed system functions.

1.5 Analysis and design methodology

1.5.1 Introduction

Large enterprise applications - the ones that execute core business applications, and keep a company going - must be more than just a bunch of code modules. They must be structured in a way that enables scalability, security, and robust execution under stressful conditions, and their structure - frequently referred to as their architecture - must be defined clearly enough that maintenance programmers can (quickly!) find and fix a bug that shows up long after the original authors have moved on to other projects. That is, these programs must be designed to work perfectly in many areas, and business functionality is not the only one (although it certainly is the essential core). Of course a well-designed architecture benefits any program, and not just the largest ones as we've singled out here. We mentioned large applications first because structure is a way of dealing with complexity, so the benefits of structure (and of modeling and design, as we'll demonstrate) compound as application size grows large. Another benefit of structure is that it enables code reuse: Design time is the easiest time to structure an application as a collection of self-contained modules or components. Eventually, enterprises build up a library of models of components, each one representing an implementation stored in a library of code modules. When another application needs the same functionality, the designer can quickly import its module from the library. At coding time, the developer can just as quickly import the code module into the application.

Modeling is the designing of software applications before coding. Modeling is an Essential Part of large software projects, and helpful to medium and even small projects as well. A model plays the analogous role in software development that blueprints and other plans (site maps, elevations, physical models) play in the building of a skyscraper. Using a model, those responsible for a software development project's success can assure themselves that business functionality is complete and correct, end-user needs are met, and program design supports requirements for scalability, robustness, security, extendibility, and other characteristics, before implementation in code renders changes difficult and expensive to make. Surveys show that large software projects have a huge probability of failure - in fact, it's more likely that a large software application will fail to meet all of its requirements on time and on budget than that it will succeed.

1.5.2 Abstraction

Models help us by letting us work at a higher level of abstraction. A model may do this by hiding or masking details, bringing out the big picture, or by focusing on different aspects of the

prototype. In UML 2.0, we can zoom out from a detailed view of an application to the environment where it executes, visualizing connections to other applications or, zoomed even further, to other sites. Alternatively, we can focus on different aspects of the application, such as the business process that it automates, or a business rules view. The new ability to nest model elements, added in UML 2.0, supports this concept directly.

The OMG's Unified Modeling Language (UML) helps specify, visualize, and document models of software systems, including their structure and design, in a way that meets all of these requirements. (We can use UML for business modeling and modeling of other non-software systems too.) Using any one of the large number of UML-based tools on the market, we can analyze our future application's requirements and design a solution that meets them, representing the results using UML 2.0's thirteen standard diagram types.

We can model just about any type of application, running on any type and combination of hardware, operating system, programming language, and network, in UML. Its flexibility let us model distributed applications that use just about any middleware on the market. Built upon fundamental OO concepts including class and operation, it's a natural fit for object-oriented languages and environments such as C++, Java, but we can use it to model non-OO applications as well in, for example, Fortran, VB, or COBOL. UML Profiles (that is, subsets of UML tailored for specific purposes) help us model Transactional, Real-time, and Fault-Tolerant systems in a natural way.

1.5.3 Models VS. Methodologies

The process of gathering and analyzing an application's requirements, and incorporating them into a program design, is a complex one and the industry currently supports many methodologies that define formal procedures specifying how to go about it. One characteristic of UML - in fact, the one that enables the widespread industry support that the language enjoys - is that it is methodology-independent. Regardless of the methodology that we use to perform our analysis and design, we can use UML to express the results. And, using XMI (XML Metadata Interchange, another OMG standard), we can transfer our UML model from one tool into a repository, or into another tool for refinement or the next step in our chosen development process. These are the benefits of standardization! UML 2.0 defines thirteen types of diagrams, divided into three categories: Six diagram types represent static application structure; three represent general types of behavior; and four represent different aspects of interactions:

Structure Diagrams include the Class Diagram, Object Diagram, Component Diagram, Composite Structure Diagram, Package Diagram, and Deployment Diagram. Behavior Diagrams include the Use Case Diagram (used by some methodologies during requirements gathering); Activity Diagram, and State Machine Diagram. Interaction Diagrams, all derived from the more general Behavior Diagram, include the Sequence Diagram, Communication Diagram, Timing Diagram, and Interaction Overview Diagram.

1.5.4 UML 2.0

OMG have already integrated the new features into this writeup, and here's a summary:

Nested Classifiers: This is an extremely powerful concept. In UML, almost every model building block we work with (classes, objects, components, behaviors such as activities and state machines, and more) is a classifier. In UML 2.0, we can nest a set of classes inside the component that manages them, or embed a behavior (such as a state machine) inside the class or component that implements it. This capability also let us build up complex behaviors from simpler ones, the capability that defines the Interaction Overview Diagram. We can layer different levels of abstraction in multiple ways: For example, we can build a model of our Enterprise, and zoom in to embedded site views, and then to departmental views within the site, and then to applications within a department. Alternatively, we can nest computational models within a business process model. **OMG's Business Enterprise Integration Domain Task Force (BEI DTF)** is currently working on several interesting new standards in business process and business rules.

Improved Behavioral Modeling: In UML 1.X, the different behavioral models were independent, but in UML 2.0, they all derive from a fundamental definition of a behavior (except for the Use Case, which is subtly different but still participates in the new organization).

Improved relationship between Structural and Behavioral Models: As we pointed out under Nested Classifiers, UML 2.0 let us designate that a behavior represented by (for example) a State Machine or Sequence Diagram is the behavior of a class or a component. That is, the new language goes well beyond the Classes and Objects well-modeled by UML 1.X to add the capability to represent not only behavioral models, but also architectural models, business process and rules, and other models used in many different parts of computing and even non-computing disciplines.

1.6 Conclusion

Two main parts are waiting to be implemented, the first one is the lambda architecture and the data analytics strategy. the second one is the microservices approach and the application logic.

CHAPTER

2

REQUIREMENTS ANALYSIS AND SPECIFICATION

After defining the project scope and making a preliminary study. In this chapter, we will deep dive in two parts, one is the functional branchy part which presents the application logic and the other one is the technical branch which present the infrastructure support.

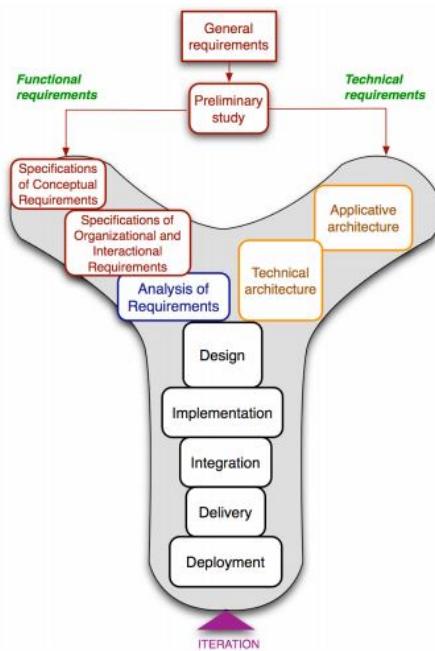


FIGURE 2.1. Sofrecom's subsidiaries and customers. Figure taken from [1].

2.1 Functional branch

The functional (left) branch corresponds to the traditional task of domain and user requirements modeling, independently from technical aspects. Considering the design of a Mixed Reality system. It includes interaction scenarios, task analysis, interaction modality choices and mock-ups. This branch ends by structuring the domains with Interactional and Business objects required to implement the Mixed Reality system.

2.1.1 Requirements analysis

2.1.1.1 Functional requirements

Our application should be able to provide the following requirements

ID	Nom	Description
1	Track users behavior	As administrator i would like to be able to track users depending on thier geolocations and interests
1	Account managment	As administrator i would like to be able to track users depending on thier geolocations and interests
1	Promotion managment	As administrator i would like to be able to track users depending on thier geolocations and interests
2	Recommend promotions and services	As administrator i would like to be able to automatically recommend services, promotions depending on users locations and interests.
3	Analyse users data flow	As administrator i would like to be able to analyse and extract informations from users data flow
4	Promote a service	As third party i would like to be able to promote my services and target a users segment
5	Benifit from services and promotions	As a client application user i would like to be able to see near promotions and services depending on my interests

2.1.1.2 Non functional requirements

Scalability Scalability is the capability of a system, network, or process to handle a growing amount of work, or its potential to be enlarged to accommodate that growth. For example, a system is considered scalable if it is capable of increasing its total output under an increased load when resources (typically hardware) are added. An analogous meaning is implied when the word is used in an economic context, where a company's scalability implies that the underlying business model offers the potential for economic growth within the company.

Operability Operability is the ability to keep an equipment, a system or a whole industrial installation in a safe and reliable functioning condition, according to pre-defined operational requirements.

Security Security is the degree of resistance to, or protection from, harm. It applies to any vulnerable and/or valuable asset, such as a person, dwelling, community, item, nation, or organization.

Fault tolerance Fault tolerance is the property that enables a system to continue operating properly in the event of the failure of (or one or more faults within) some of its components. If its operating quality decreases at all, the decrease is proportional to the severity of the failure, as compared to a naively designed system in which even a small failure can cause total breakdown. Fault tolerance is particularly sought after in high-availability or life-critical systems. The ability of maintaining functionality when portions of a system break down is referred to as graceful degradation.

Extensibility Extensibility is a software engineering and systems design principle where the implementation takes future growth into consideration. The term extensibility can also be seen as a systemic measure of the ability to extend a system and the level of effort required to implement the extension. Extensions can be through the addition of new functionality or through modification of existing functionality.

2.1.2 Requirements specification

In order to properly establish the planned objectives during our development process, we will specify in more detail the previous requirements using the UML design methodology.

2.1.2.1 Actors identification

Starting from the last requirements collection , we identify three main actors that interact continuously with our application.

- Administrator - the administrator is able to see users data analytics , track them spatially and manage users , promotions and services .

- Third party - the third party is able to submit a promotion via an application portal.
- User - the user is able to see promotions and services depending on his interests.

2.1.2.2 Use case modelization

This phase aims to describe the expected behavior of the application. For this purpose, we use the use case diagram as an essential element of object-oriented modeling. It makes it possible to model the functionalities of the application .

Data analytics use case diagram :

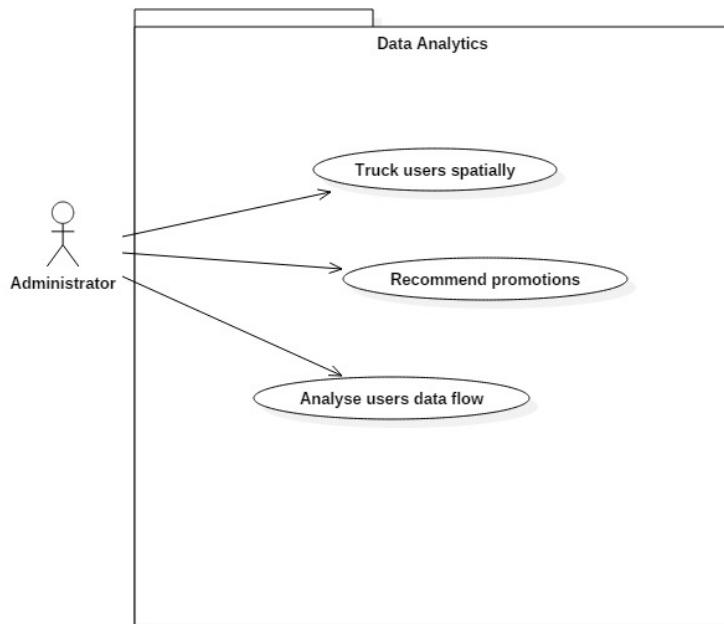


FIGURE 2.2. Sofrecom's subsidiaries and customers. Figure taken from [1].

Use case identification :

Title : Data analytics

Summary : Our application administrator is able to truck users spatially and analyse their data flow , relaying on the big data engine , to recommend the most apporiate promotion depending on thier intersts and locations.

Actors : administrator

Scenarios description :

Precondition : This functionality user require the administrator privilege.

Post-condition : Users could receive administrator logic promotions recommendation.

Nominal Scenario :

1. Administrator login to dashboard
2. Administrator choose between Dashboard and World Map
3. If Administrator choose dashboard
 - a) administrator is able to see data flow analytics in deferent type of charts
4. If Administrator choose world map
 - a) administrator is able to see users Heatmap

Alternative Scenario :

There is no alternative scenario.

Error Scenario :

If the dashboard could not get informations from the microservice then the dashboard show initial information.

Use case diagram of account management :

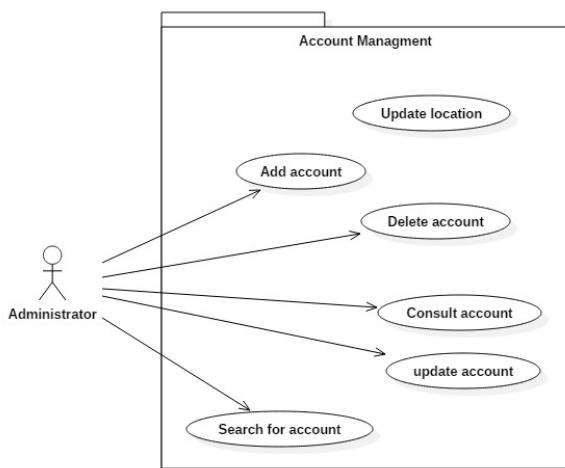


FIGURE 2.3. Sofrecom's subsidiaries and customers. Figure taken from [1].

Use case identification :

Title : Account management

Summary : Our application administrator is able to add, remove, update , search and consult users. Our system provide an update position functionality that allows real time user position tracking. Mobile client application send periodically the device position to the account microservice.

Actors : administrator

Scenarios description :

Precondition : This functionality user require the administrator privilege.

Post-condition : Administrator could mange users account and track their positions.

Nominal Scenario :

1. Administrator login to dashboard
2. Administrator choose account management
3. Administrator could search, add, update, delete and consult an account.

Alternative Scenario : there is no alternative scenario.

Error Scenario : If the administrator was not able to do the opration then he will be notified.

Promotions managment :

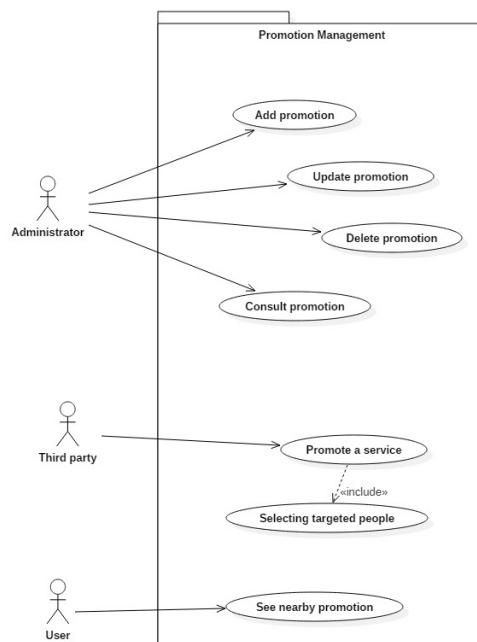


FIGURE 2.4. Sofrecom's subsidiaries and customers. Figure taken from [1].

Use case identification :

Title : Promotion management

Summary : Our application administrator is able to search, add, delete, update and consult promotions, Third party representative is able to make to promote a serive and select which kind of people his intersted about and user can see promotion and take advantage of promotions services.

Actors : administrator , Third party and user

Scenarios description :

Precondition : This functionality user require a varity of privilege.

Post-condition : Promotions are well exposed to app users.

Nominal Scenario :

1. Administrator login to dashboard
2. Administrator choose Promotion management
3. administrator is able to search, add, update, delete, consult promotion
1. Third party representative login to dashboard
2. Choose Promotion management
3. He is able to search, add, update, delete and consult his promotions
1. User login to the mobile application
2. User is able to see promotions

Alternative Scenario : There is no alternative senario.

Error Scenario : If the any application user (administrator, third party representative, mobile application user) was not able to do the opration then he will be notified.

2.1.2.3 Sequence diagrams

With sequence diagrams, the UML provides a graphical way to represent interactions between objects over time. These diagrams typically show an actor, the objects and the actors with whom he interacts during the execution of the use case. In this section, we present some sequence diagrams to describe the different interactions between the user and the application

Sequence Diagram for the Authentication Scenario :

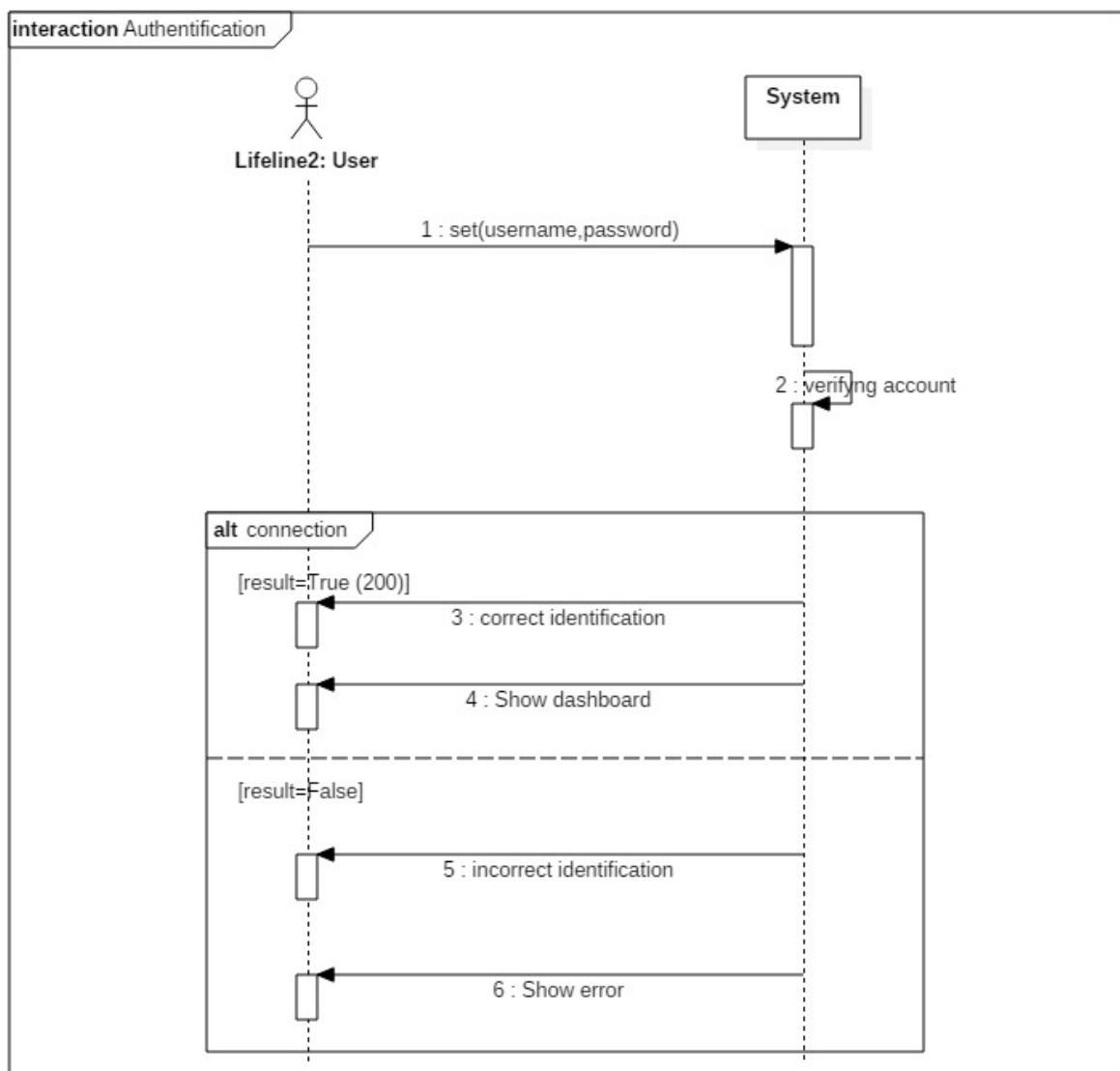


FIGURE 2.5. Sofrecom's subsidiaries and customers. Figure taken from [1].

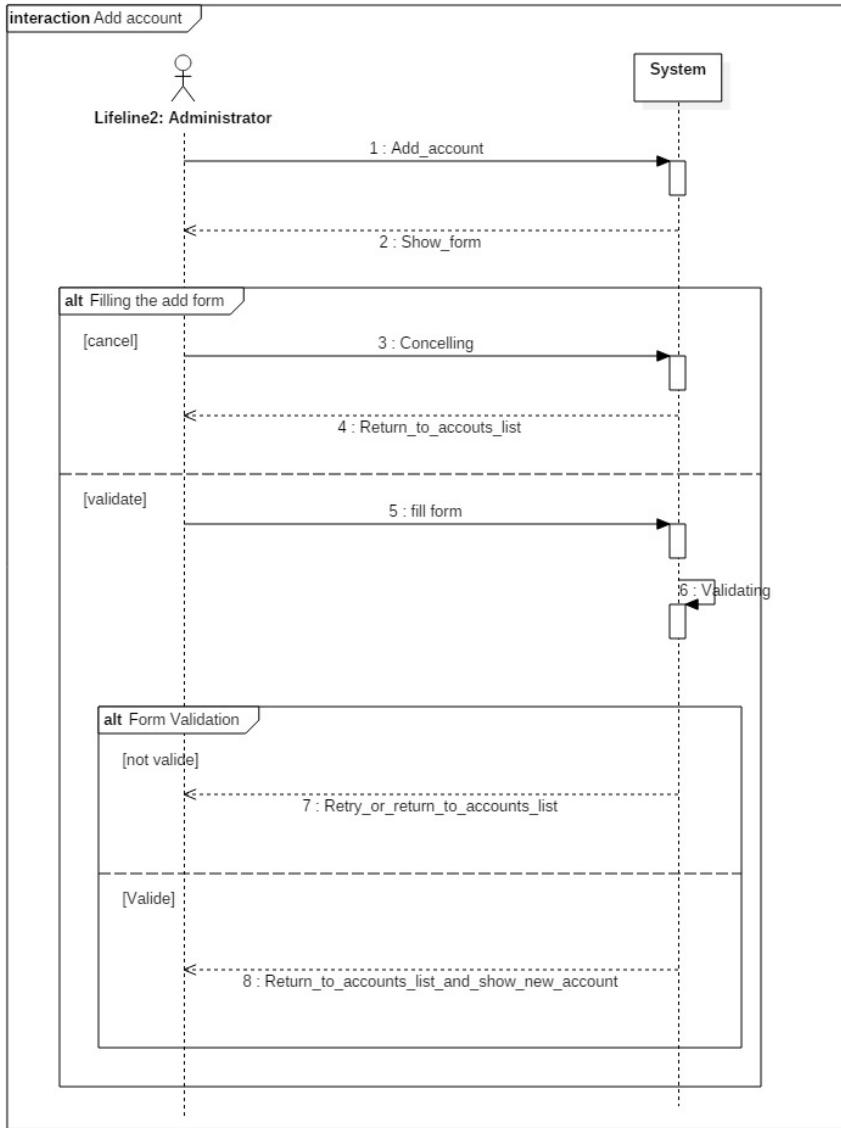
Sequence Diagram for account management :


FIGURE 2.6. Sofrecom's subsidiaries and customers. Figure taken from [1].

Conclusion this phase allows us to initiate to the following chapter which deals with our application design. In the next section we will deep dive in technical specification.

2.2 Technical branch

Since the functional requirements have been defined, it's time to choose which tools shall we use and especially on which software architecture we will be reposing.

2.2.1 Big Data architecture

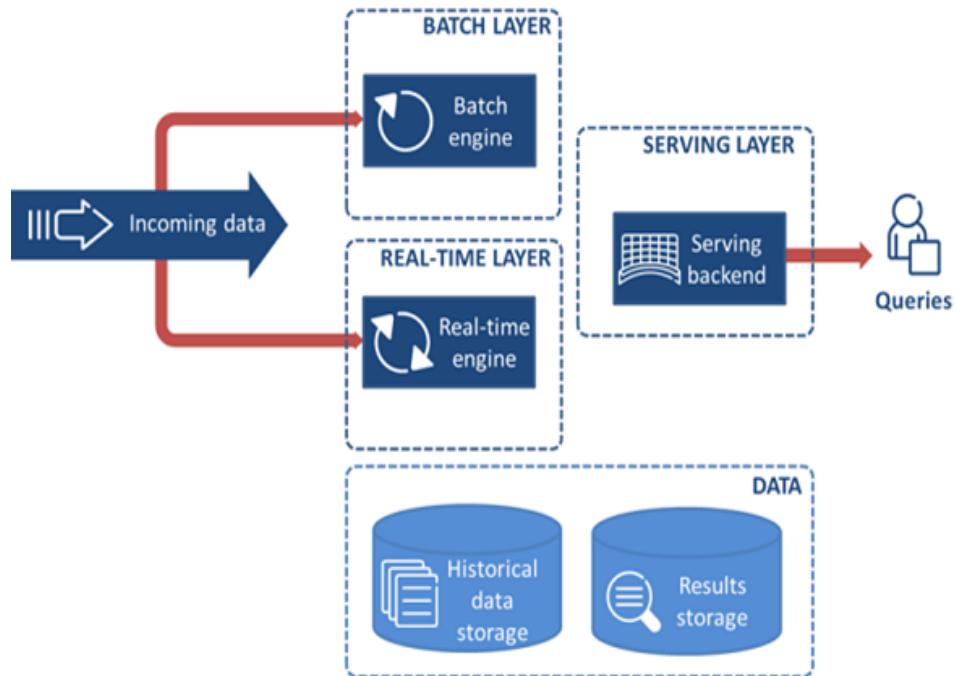


FIGURE 2.7. Sofrecom's subsidiaries and customers. Figure taken from [1].

In the last chapter, we choosed to go with the lambda approach. the question now is how to choose the right tools. So in the next sections we will be making a comparative study in each field (yellow tags on the previous figure) to pick the most appropriate tools.

2.2.1.1 OSS Message brokers layer

The incoming data could be from different sources . Data should be gathered in one data flow pipeline. Otherwise we will find ourselves dealing with a complexe grid , providing all apis compatibility and resolving extras problems.

A great solution for this problem could be a messaging layer that gives a various guarantees of message persistence and delivery.In this section , we will introduce some of the best message brokers that can handle this activity.So our Competitors are SQS AWS, Apache Kafka, Mongodb message queue, HornetQ, RabbitMQ, Apache ActiveMQ

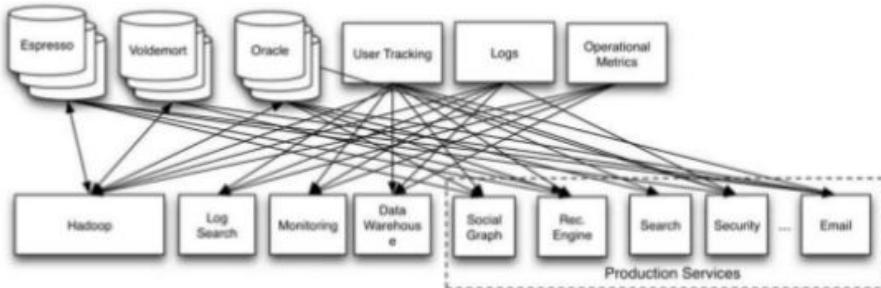


FIGURE 2.8. Sofrecom's subsidiaries and customers. Figure taken from [1].

SQS AWS SQS, Simple Message Queue, is a message-queue-as-a-service offering from Amazon Web Services. It supports only a handful of messaging operations but thanks to the easy to understand interfaces, and the as-a-service nature. SQS guarantees that if a send completes, the message is replicated to multiple nodes. It also provides at-least-once delivery guarantees. We don't really know how SQS is implemented, A single thread on single node achieves 430 msgs/s sent and the same number of msgs received. These results are not impressive, but SQS scales nicely both when increasing the number of threads, and the number of nodes . On a single node, with 50 threads, we can send up to 14 500 msgs/s, and receive up to 4 200 msgs/s. On an 8-node cluster, these numbers go up to 63 500 msgs/s sent, and 34 800 msgs/s received.

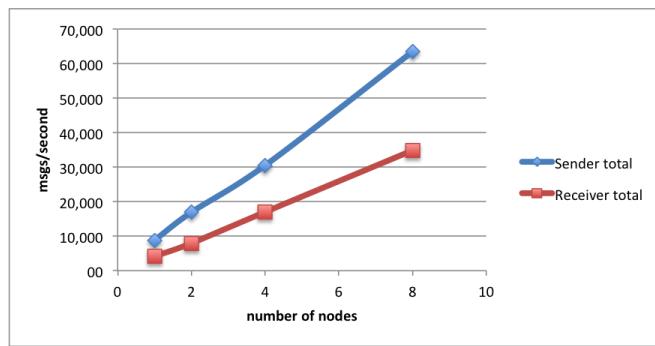


FIGURE 2.9. Sofrecom's subsidiaries and customers. Figure taken from [1].

Apache Kafka Apache Kafka is an open-source stream processing platform developed by the Apache Software Foundation written in Scala and Java. The project aims to provide a unified, high-throughput, low-latency platform for handling real-time data feeds. Its storage layer is essentially a massively scalable pub/sub message queue architected making it highly valuable for enterprise infrastructures to process streaming data. Kafka's performance is great. A single-node single-thread achieves about 2 550 msgs/s, and the best result was 33 500msgs/s with 25 sending and receiving threads and 4 nodes.

THREADS	NODES	SEND MSGS/S	RECEIVE MSGS/S
1	1	2 558,81	2 561,10
5	1	11 804,39	11 354,90
25	1	29 691,33	27 093,58
1	2	5 879,21	5 847,03
5	2	14 660,90	13 552,35
25	2	29 373,50	27 822,50
5	4	22 271,50	20 757,50
25	4	33 587,00	31 891,00

FIGURE 2.10. Sofrecom's subsidiaries and customers. Figure taken from [1].

Mongodb message queue Mongo has two main features which make it possible to easily implement a durable, replicated message queue on top of it. A single-thread, single-node setup achieves 7 900 msgs/s sent and 1 900 msgs/s received. The maximum send throughput with multiple thread/nodes is about 10 500 msgs/s, while the maximum receive rate is 3 200 msgs/s, when using the „safe” write concern.

THREADS	NODES	SEND MSGS/S	RECEIVE MSGS/S
1	1	7 968,60	1 914,05
5	1	9 903,47	3 149,00
25	1	10 903,00	3 266,83
1	2	9 569,99	2 779,87
5	2	10 078,65	3 112,55
25	2	7 930,50	3 014,00

FIGURE 2.11. Sofrecom's subsidiaries and customers. Figure taken from [1].

HornetQ HornetQ is an open source asynchronous messaging project from JBoss. It is an example of Message-oriented middleware. HornetQ is an open source project to build a multi-protocol, embeddable, very high performance, clustered, asynchronous messaging system. Speaking of performance, HornetQ is very good! A single-node, single-thread setup achieves 1 100 msgs/s. With 25 threads, we are up to 12 800 msgs/s! And finally, with 25 threads and 4 nodes, we can achieve 17 000 msgs/s.

2.2. TECHNICAL BRANCH

THREADS	NODES	SEND MSGS/S	RECEIVE MSGS/S
1	1	1 108,38	1 106,68
5	1	4 333,13	4 318,25
25	1	12 791,83	12 802,42
1	2	2 095,15	2 029,99
5	2	7 855,75	7 759,40
25	2	14 200,25	13 761,75
1	4	3 768,28	3 627,02
5	4	11 572,10	10 708,70
25	4	17 402,50	16 160,50

FIGURE 2.12. Sofrecom's subsidiaries and customers. Figure taken from [1].

RabbitMQ RabbitMQ is one of the leading open-source messaging systems. It is written in Erlang, implements AMQP and is a very popular choice when messaging is involved. It supports both message persistence and replication, with well documented behaviour in case of e.g. partitions. Such strong guarantees are probably one of the reasons for poor performance. A single-thread, single-node gives us 310 msgs/s sent and received. This scales nicely as we add nodes, up to 1 600 msgs/s

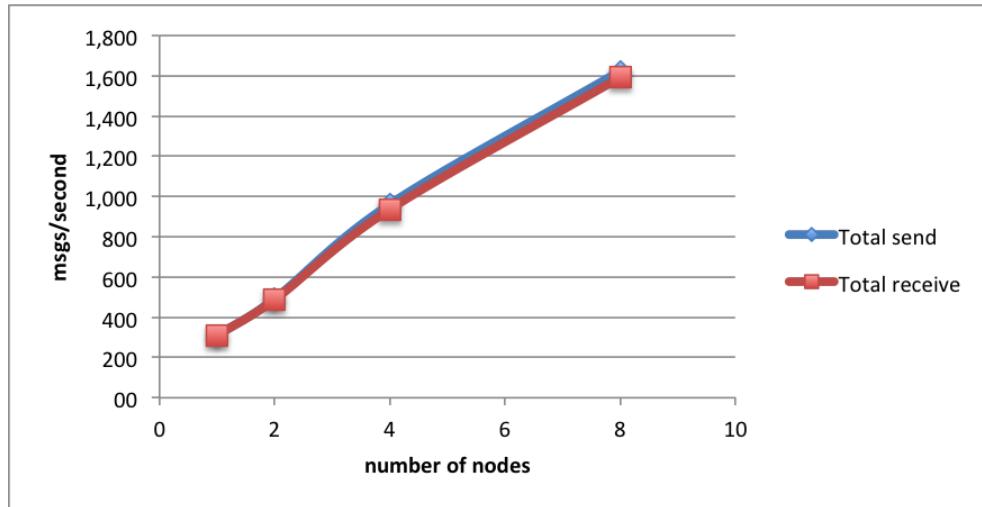


FIGURE 2.13. Sofrecom's subsidiaries and customers. Figure taken from [1].

Apache ActiveMQ Apache ActiveMQ is an open source message broker written in Java together with a full Java Message Service (JMS) client. It provides Enterprise Features which in

this case means fostering the communication from more than one client or server. Supported clients include Java via JMS 1.1 as well as several other "cross language" clients. The communication is managed with features such as computer clustering and ability to use any database as a JMS persistence provider besides virtual memory, cache, and journal persistency. The next benchmark show the performance of ActiveMQ :

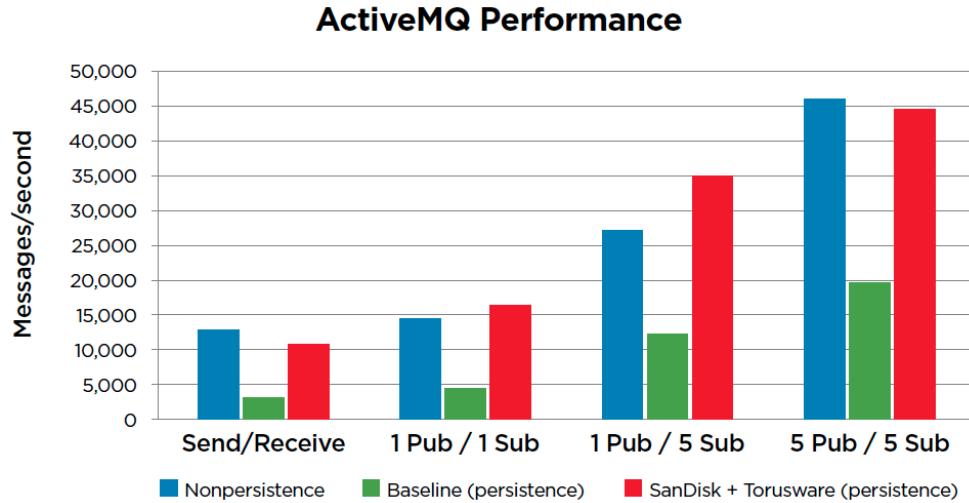


FIGURE 2.14. Sofrecom's subsidiaries and customers. Figure taken from [1].

Summary Now the question is which message queue we should choose . All of the above solutions have some good sides:

- SQS is a service, so especially if you are using the AWS cloud, it's an easy choice: good performance and no setup required
- if you are using Mongo, it is easy to build a replicated message queue on top of it, without the need to create and maintain a separate messaging cluster
- if you want to have high persistence guarantees, RabbitMQ ensures replication across the cluster and on disk on message send.
- HornetQ has great performance with a very rich messaging interface and routing options
- Kafka offers the best performance and scalability
- ActiveMQ employs several modes for high availability,A robust horizontal scaling mechanism and a flexibility in configuration

When looking only at the throughput, Kafka is a clear winner (unless we include SQS with multiple nodes, but as mentioned, that would be unfair)

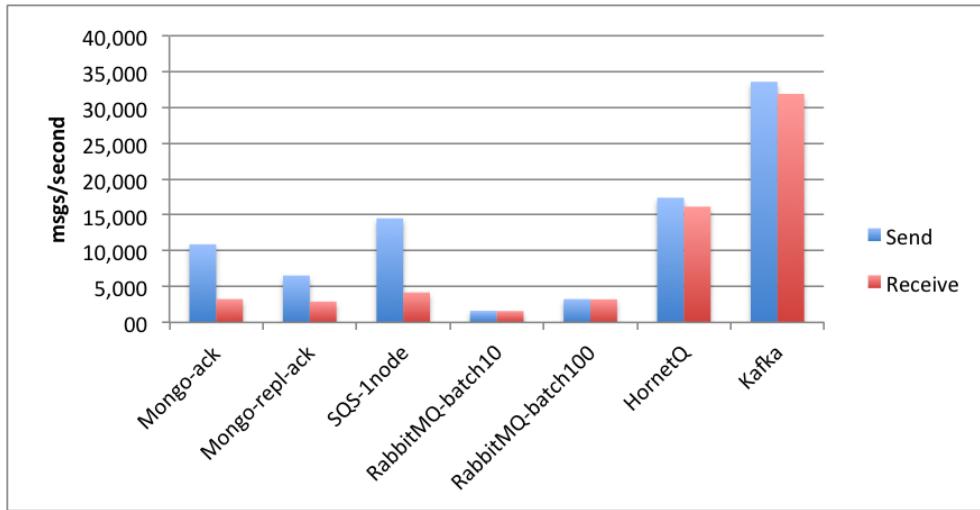


FIGURE 2.15. Sofrecom's subsidiaries and customers. Figure taken from [1].

It's also interesting to see how sending more messages in a batch improves the throughput, when increasing the batch size from 10 to 100, Rabbit gets a 2x speedup, HornetQ a 1.2x speedup, and Kafka a 2.5x speedup, achieving about 89 000 msgs/s!

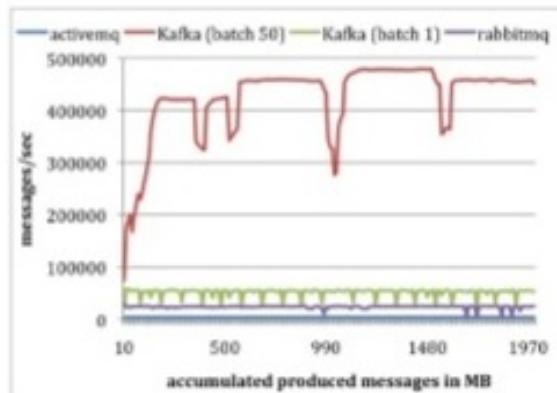


Figure 4. Producer Performance

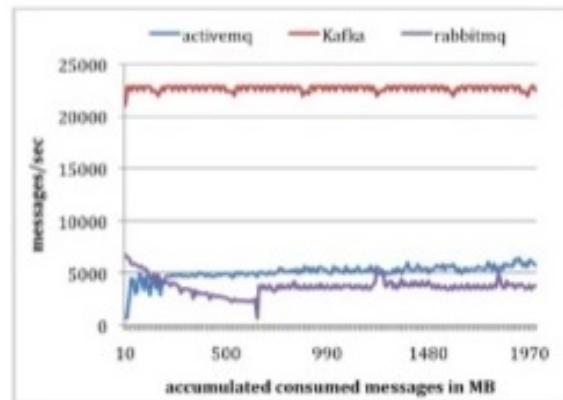


Figure 5. Consumer Performance

FIGURE 2.16. Sofrecom's subsidiaries and customers. Figure taken from [1].

There are of course many other aspects besides performance, which should be taken into account when choosing a message queue, such as administration overhead, partition tolerance, feature set regarding routing, etc. And as we are interested in first with performance , we will choose Kafka for our mission .

2.2.1.2 NoSQL Databases

After performing the Data analytics logic , results are persisted in a serving layer. This layer is mostly a databases and more precisely a noSql database. In this section , we will introduce some of the best noSql databases that can handle this activity. So our Competitors are:

- **Document oriented :** Mongodb , CouchDB
- **Key Value oriented :** Redis
- **Graph oriented :** Neo4j
- **Column oriented :** Cassandra , Riak , HBase

Document oriented :

A document-oriented database, or document store, is a computer program designed for storing, retrieving and managing document-oriented information, also known as semi-structured data. Document-oriented databases are one of the main categories of NoSQL databases, and the popularity of the term document-oriented database has grown with the use of the term NoSQL itself.

Mongodb Mongodb is one of the most popular document based NoSQL database as it stores data in JSON like documents. It is non-relational database with dynamic schema. It has been developed by the founders of DoubleClick, written in C++ and is currently being used by some big companies like The New York Times, Craigslist, MTV Networks. The following are some of MongoDB benefits and strengths :

- Speed: For simple queries, it gives good performance, as all the related data are in single document which eliminates the join operations.
- Scalability: It is horizontally scalable i.e. you can reduce the workload by increasing the number of servers in your resource pool instead of relying on a stand alone resource.
- Manageable: It is easy to use for both developers and administrators. This also gives the ability to shard database
- Dynamic Schema: Its gives you the flexibility to evolve your data schema without modifying the existing data

CouchDB CouchDB is also a document based NoSQL database. It stores data in form of JSON documents. The following are some of CouchDB benefits and strengths:

- Schema-less: As a member of NoSQL family, it also have dynamic schema which makes it more flexible, having a form of JSON documents for storing data.

- HTTP query: You can access your database documents using your web browser.
- Conflict Resolution: It has automatic conflict detection which is useful while in a distributed database.
- Easy Replication: Implementing replication is fairly straight forward

Key Value oriented :

A key-value store, or key-value database, is a data storage paradigm designed for storing, retrieving, and managing associative arrays, a data structure more commonly known today as a dictionary or hash. Dictionaries contain a collection of objects, or records, which in turn have many different fields within them, each containing data. These records are stored and retrieved using a key that uniquely identifies the record, and is used to quickly find the data within the database.

Redis Redis is an Open Source NoSQL database which is mainly used because of its lightening speed. It is written in ANSI C language. The following are some of Redis benefits and strengths:

- Data structures: Redis provides efficient data structures to an extend that it is sometimes called as data structure server. The keys stored in database can be hashes, lists, strings, sorted or unsorted sets.
- Redis as Cache: You can use Redis as a cache by implementing keys with limited time to live to improve the performance.
- Very fast: It is consider as one of the fastest NoSQL server as it works with the in-memory dataset.

Graph oriented :

The graph based DBMS models represent the data in a completely different way than the previous three models. They use tree-like structureswith nodes and edges connecting each other through relations.

Neo4j Neo4j is a graph database management system developed by Neo Technology,Described by its developers as an ACID-compliant transactional database with native graph storage and processing, Neo4j is the most popular graph database. Neo4j is implemented in Java and accessible from software written in other languages using the Cypher Query Language through a transactional HTTP endpoint, or through the binary bolt protocol.

Column oriented :

A column-oriented DBMS (or columnar database management system) is a database management system (DBMS) that stores data tables by column rather than by row. Practical use of a column store versus a row store differs little in the relational DBMS world. Both columnar and row databases can use traditional database query languages like SQL to load data and perform queries.

Cassandra Apache Cassandra is the leading NoSQL, distributed database management system driving many of today's modern business applications by offering continuous availability, high scalability and performance, strong security, and operational simplicity while lowering overall cost of ownership. Cassandra has decentralized architecture. Any node can perform any operation. It provides AP(Availability,Partition-Tolerance). Cassandra has excellent single-row read performance as long as eventual consistency semantics are sufficient for the use-case.

Riak Riak is a distributed NoSQL key-value data store that offers high availability, fault tolerance, operational simplicity, and scalability. In addition to the open-source version, it comes in a supported enterprise version and a cloud storage version. Riak implements the principles from Amazon's Dynamo paper with heavy influence from the CAP Theorem. Written in Erlang, Riak has fault tolerance data replication and automatic data distribution across the cluster for performance and resilience.

HBase HBase is an open source, non-relational, distributed database modeled after Google's Bigtable and is written in Java. It is developed as part of Apache Software Foundation's Apache Hadoop project and runs on top of HDFS (Hadoop Distributed File System), providing Bigtable-like capabilities for Hadoop. That is, it provides a fault-tolerant way of storing large quantities of sparse data. HBase is a column-oriented key-value data store and has been idolized widely because of its lineage with Hadoop and HDFS. HBase runs on top of HDFS and is well-suited for faster read and write operations on large datasets with high throughput and low input/output latency.

Summary :

It's too hard to pick the right database for our serving layer . In this paragraph we try to choose the most appropriate one depending on some benchmarks well known. This benchmarks elect

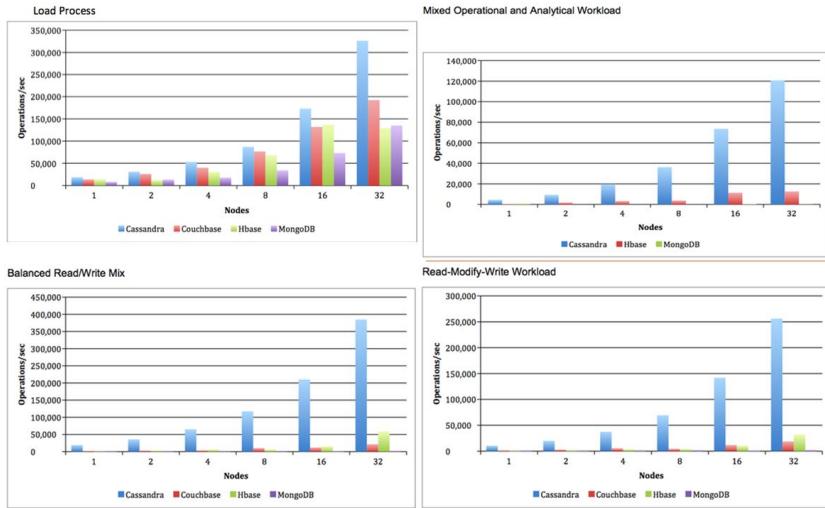


FIGURE 2.17. Sofrecom's subsidiaries and customers. Figure taken from [1].

Cassandra as the best database in performance point of view. But in fact , there is no best database , choosing the best one depend on the use case . It's the Polyglot persistence which we'll explain in the next section.

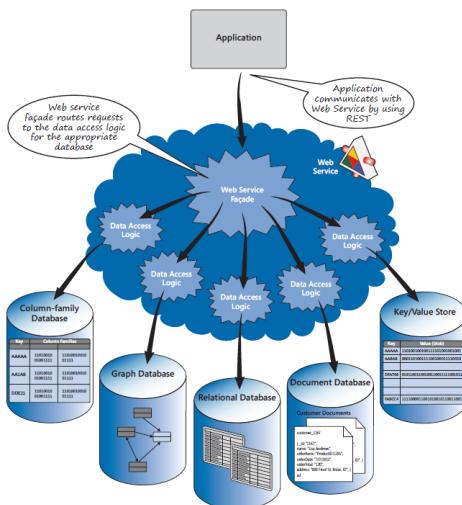
Polyglot persistence :


FIGURE 2.18. Sofrecom's subsidiaries and customers. Figure taken from [1].

Polyglot Persistence is a fancy term to mean that when storing data, it is best to use multiple data storage technologies, chosen based upon the way data is being used by individual applications or components of a single application. Different kinds of data are best dealt with different data stores. In short, it means picking the right tool for the right use case. It's the same idea behind Polyglot Programming, which is the idea that applications should be written in a mix of languages to take advantage of the fact that different languages are suitable for tackling different problems.



FIGURE 2.19. Sofrecom's subsidiaries and customers. Figure taken from [1].

And as we want to make data analytics. We choose as a noSQL database for our system , the oriented column one "Cassandra".

2.2.1.3 Big Data Processing Frameworks

Only-Batch processing frameworks

Apache Hadoop (MapReduce) Modern versions of Hadoop are composed of several components or layers, that work together to process batch data:

- **HDFS:** HDFS is the distributed filesystem layer that coordinates storage and replication across the cluster nodes. HDFS ensures that data remains available in spite of inevitable host failures. It is used as the source of data, to store intermediate processing results, and to persist the final calculated results.
- **YARN:** YARN, which stands for Yet Another Resource Negotiator, is the cluster coordinating component of the Hadoop stack. It is responsible for coordinating and managing the underlying resources and scheduling jobs to be run. YARN makes it possible to run much more diverse workloads on a Hadoop cluster than was possible in earlier iterations by acting as an interface to the cluster resources.

- MapReduce: MapReduce is Hadoop's native batch processing engine.
- Hadoop commons : librerie that supports others component.

MapReduce's processing technique follows the map, shuffle, reduce algorithm using key-value pairs.

Apache Hadoop and its MapReduce processing engine offer a well-tested batch processing model that is best suited for handling very large data sets where time is not a significant factor.

Only-Streaming processing frameworks

Apache storm Apache Storm is a stream processing framework that focuses on extremely low latency and is perhaps the best option for workloads that require near real-time processing. It can handle very large quantities of data with and deliver results with less latency than other solutions. Storm stream processing works by orchestrating DAGs (Directed Acyclic Graphs) in a framework it calls topologies. The topologies are composed of:

- Streams: Conventional data streams. This is unbounded data that is continuously arriving at the system.
- Spouts: Sources of data streams at the edge of the topology. These can be APIs, queues, etc. that produce data to be operated on.
- Bolts: Bolts represent a processing step that consumes streams, applies an operation to them, and outputs the result as a stream. Bolts are connected to each of the spouts, and then connect to each other to arrange all of the necessary processing. At the end of the topology, final bolt output may be used as an input for a connected system.

Storm is probably the best solution currently available for near real-time processing. It is able to handle data with extremely low latency for workloads that must be processed with minimal delay. Storm is often a good choice when processing time directly affects user experience, for example when feedback from the processing is fed directly back to a visitor's page on a website.

Storm with Trident gives you the option to use micro-batches instead of pure stream processing. While this gives users greater flexibility to shape the tool to an intended use, it also tends to negate some of the software's biggest advantages over other solutions. That being said, having a choice for the stream processing style is still helpful. For pure stream processing workloads with very strict latency requirements, Storm is probably the best mature option. It can guarantee message processing and can be used with a large number of programming languages. Because Storm does not do batch processing, you will have to use additional software if you require those capabilities. If you have a strong need for exactly-once processing guarantees, Trident can provide that. However, other stream processing frameworks might also be a better fit at that point.

Samza Apache Samza is a good choice for streaming workloads where Hadoop and Kafka are either already available or sensible to implement. Samza itself is a good fit for organizations with multiple teams using (but not necessarily tightly coordinating around) data streams at various stages of processing. Samza greatly simplifies many parts of stream processing and offers low latency performance. It might not be a good fit if the deployment requirements aren't compatible with your current system, if you need extremely low latency processing, or if you have strong needs for exactly-once semantics.

Hybrid processing frameworks

Apache Spark Spark is a great option for those with diverse processing workloads. Spark batch processing offers incredible speed advantages, trading off high memory usage. Spark Streaming is a good stream processing solution for workloads that value throughput over latency.

Apache Flink Flink offers both low latency stream processing with support for traditional batch tasks. Flink is probably best suited for organizations that have heavy stream processing requirements and some batch-oriented tasks. Its compatibility with native Storm and Hadoop programs, and its ability to run on a YARN-managed cluster can make it easy to evaluate. Its rapid development makes it worth keeping an eye on.

Summary :



Streaming Model	Native	Micro-batching	Micro-batching	Native	Native
API	Compositional		Declarative	Compositional	Declarative
Guarantees	At-least-once	Exactly-once	Exactly-once	At-least-once	Exactly-once
Fault Tolerance	Record ACKs		RDD based Checkpointing	Log-based	Checkpointing
State Management	Not build-in	Dedicated Operators	Dedicated DStream	Stateful Operators	Stateful Operators
Latency	Very Low	Medium	Medium	Low	Low
Throughput	Low	Medium	High	High	High
Maturity	High		High	Medium	Low

FIGURE 2.20. Sofrecom's subsidiaries and customers. Figure taken from [1].

Infact, as before there is no best processing framework . Each one of them is good for a specific usage. So three cases occur:

- Only Batch processing : Hadoop and MapReduce would be a great solution.
- Only Streaming processing : Apache Strom is now undefeatable.
- Hybrid processing : Spark would be a very good solution..

And as we would like to implement the lambda architecture and we definitely want to go with the hybrid processing option. We choose as a data processing engine for our batch and streaming layers "Apache Spark".

2.2.1.4 Conclusion

Finally, we are now able to gather all pieces into one. This comparative study supports our choices in each layer. In the next section we will talk about our microservices architecture.

2.2.2 Architecture microservice

Our application is divided in two parts, the first one is the big data part explained in the last section and the second one is the microservices part. The microservices architecture presents so many patterns. In this section, we will introduce some microservices patterns, then we choose for our application the most appropriate one. After fixing all required patterns we finally move to define our microservices architecture.

2.2.2.1 Decomposition

The microservice architecture structures an application as a set of loosely coupled services. The application should be decomposed in a way so that most new and changed requirements only affect a single service. So the question now is how to decompose an application into services? the next two paragraphs show some useful decomposition strategy .

Decompose by business capability Define services corresponding to business capabilities. A business capability is a concept from business architecture modeling. It is something that a business does in order to generate value. A business capability often corresponds to a business object.

Decompose by subdomain Define services corresponding to Domain-Driven Design (DDD) subdomains. DDD refers to the application's problem space - the business - as the domain. A domain is consists of multiple subdomains. Each subdomain corresponds to a different part of the business.

Summary Dividing services upon business capabilities is more suitable for us. The decomposition by business capability present so much benefits , some of them are :

- Stable architecture since the business capabilities are relatively stable
- Development teams are cross-functional, autonomous, and organized around delivering business value rather than technical features
- Services are cohesive and loosely coupled

2.2.2.2 Cross cutting concerns

The cross-cutting concern include:

- Externalized configuration - includes credentials, and network locations of external services such as databases and message brokers
- Logging - configuring of a logging framework such as log4j or logback
- Health checks - a url that a monitoring service can ping to determine the health of the application
- Metrics - measurements that provide insight into what the application is doing and how it is performing
- Distributed tracing - instrument services with code that assigns each external request an unique identifier that is passed between services.

There are so much services. we will frequently create new services, each of them will only take days or weeks to develop. we cannot afford to spend a few days configuring the mechanisms to handle cross-cutting concerns. What is even worse is that in a microservice architecture there are additional cross-cutting concerns that we have to deal with including service registration and discovery. So the solution is to build our microservices using a microservice chassis framework, which handles cross-cutting concerns.

Microservice chassis The major benefit of a microservice chassis is that we can quickly and easy get started with developing a microservice. we need a microservice chassis for each programming language/framework that we want to use. This can be an obstacle to adopting a new programming language or framework. There are so much microservice chassis :

- Spring Boot and Spring Cloud (java)
- Dropwizard (java)
- Gizmo (go)

- Micro (go)
- Go kit (go)

Summary We didn't make a comparative study on chassis frameworks, Spring boot and Spring Cloud have a big community and present so much quick start features. So we choose to go with Spring boot and Spring cloud as a chassis framework.

2.2.2.3 External API

Data in the Microservice architecture pattern is spread over multiple service. Each client needs to fetch information from all of these services. The question now is how do the clients of a Microservices-based application access the individual services? Exposing an Api gateway could be a great solution :

API gateway :

Implement an API gateway that is the single entry point for all clients. The API gateway handles

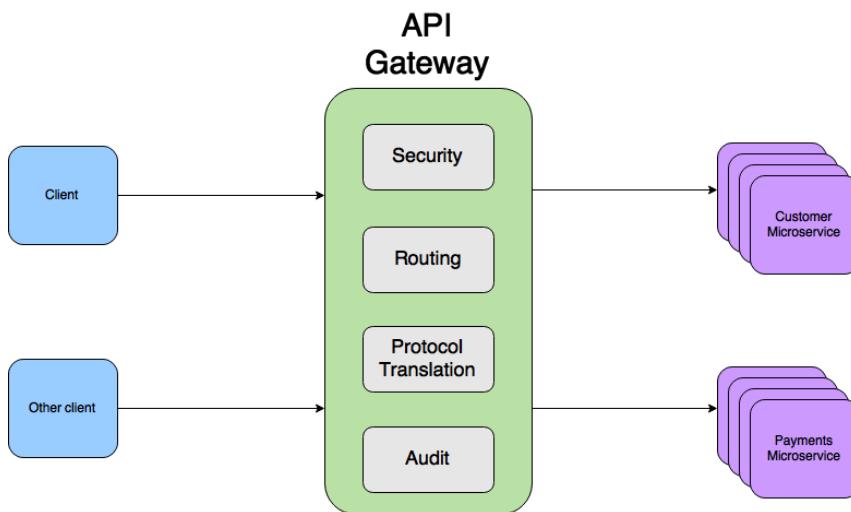


FIGURE 2.21. Sofrecom's subsidiaries and customers. Figure taken from [1].

requests in one of two ways. Some requests are simply proxied/routed to the appropriate service. It handles other requests by fanning out to multiple services. the API gateway can expose a different API for each client. The API gateway might also implement security, e.g. verify that the client is authorized to perform the request

Summary Exposing an Api gateway is a bright idea. To implement this solution there is some tools that can handle this :

- Kong API Api Gatway

- Zuul and Spring Cloud

As we choosed Spring boot and Spring Cloud for our application chassis framework, we definitely go with the "Zuul and Spring cloud" api gatway option.

2.2.2.4 Data management

Most services need to persist data in some kind of database. The question is what's the database architecture in a microservices application ?

Database per Service Keep each microservice's persistent data private to that service and accessible only via its API. The service's database is effectively part of the implementation of that service. It cannot be accessed directly by other services. There are a few different ways to keep a service's persistent data private.

- Private-tables-per-service
- Schema-per-service
- Database-server-per-service

Shared database Use a (single) database that is shared by multiple services. Each service freely accesses data owned by other services using local ACID transactions.

Summary choosing a the option of database per service could be a great solution. This pattern present some benefits :

- Helps ensure that the services are loosely coupled. Changes to one service's database does not impact any other services.
- Each service can use the type of database that is best suited to its needs

so we choose to go with the the database per service choice.

2.2.2.5 Conclusion

In this section we have fixed all patterns and choosed which technologies we would like to implement , in the next section we move to talk about our deployment strategy.

2.2.3 Deployment architecture

Our system now is a set of services. Each service is deployed as a set of service instances for throughput and availability. The question now is how are services packaged and deployed?

Multiple service instances per host :

Run multiple instances of different services on a host (Physical or Virtual machine). There are various ways of deploying a service instance on a shared host including:

- Deploy each service instance as a JVM process. For example, a Tomcat or Jetty instances per service instance.
- Deploy multiple service instances in the same JVM. For example, as web applications.

Service instance per host :

Deploy each single service instance on its own host

Service instance per VM :

Package the service as a virtual machine image and deploy each service instance as a separate VM

Service instance per Container :

Package the service as a (Docker) container image and deploy each service instance as a container

Summary :

Deploying our services in containers is the most appropriate option for us. Due resource and budget limits, our application is deployed under a VM that conatin many containers carrying services (service per container). This approach has some benefits :

- It is straightforward to scale up and down a service by changing the number of container instances.
- The container encapsulates the details of the technology used to build the service. All services are, for example, started and stopped in exactly the same way.
- Each service instance is isolated.

- Containers are extremely fast to build and start. For example, it's 100x faster to package an application as a Docker container than it is to package it as an AMI. Docker containers also start much faster than a VM since only the application process starts rather than an entire OS.
- A container imposes limits on the CPU and memory consumed by a service instance.

2.3 Conclusion

Two main parts are waiting to be implemented, the first one is the lambda architecture. The second one is the microservices approach and the application logic. In our scope , due to time limits we will not be able to implement all patterns , we will just make the most importants one. In the next chapter we move to our application Design.

Notice :

In our scope , due to time limits we will not be able to implement all functionalities and patterns , we will just make the most importants one.

CHAPTER



3

DESIGN

This chapter will be devoted to present the design stage of our system. Designing is a creative process that describes the manner in which our system works. We will start by presenting our architecture design, then we move to the general conception and finally the detailed design.

3.1 General design

Our project is essentially based on two main parts as described before. In the next section, we present the architecture of the system.

3.1.1 Architectural design

Architectural design refers to the high level structures of a software system, the discipline of creating such structures, and the documentation of these structures. These structures are needed to reason about the software system. Each structure comprises software elements, relations among them, and properties of both elements and relations. It's about making fundamental structural choices which are costly to change once implemented. Architecture design choices include specific structural options from possibilities in the design of software.

3.1.1.1 Big Data architectural design

In the last chapter we fixed what we will be using in term of patterns and tools. Now, we are able to just start implementing our lambda architecture.

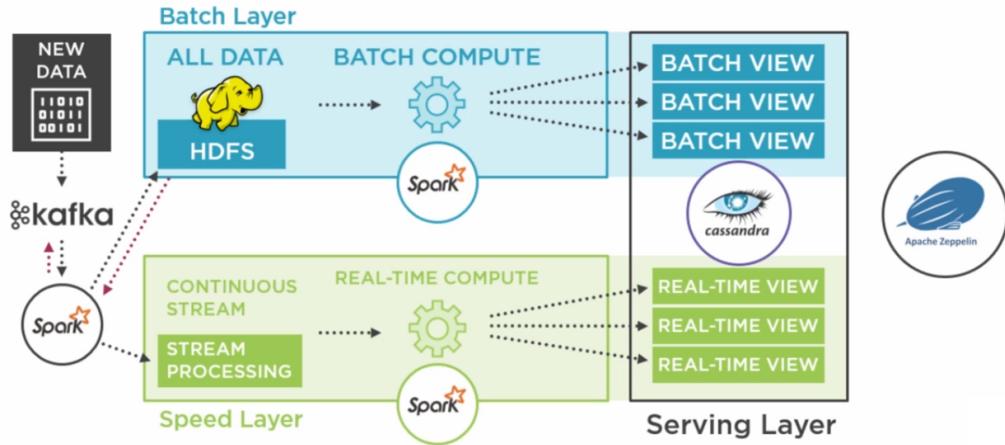


FIGURE 3.1. Sofrecom's subsidiaries and customers. Figure taken from [1].

3.1.1.2 Microservices architectural design

After fixing all patterns and choosing which technologies we would like to implement , it's time to make the last step. The next figure show our microservice application architecture in it's first versions:

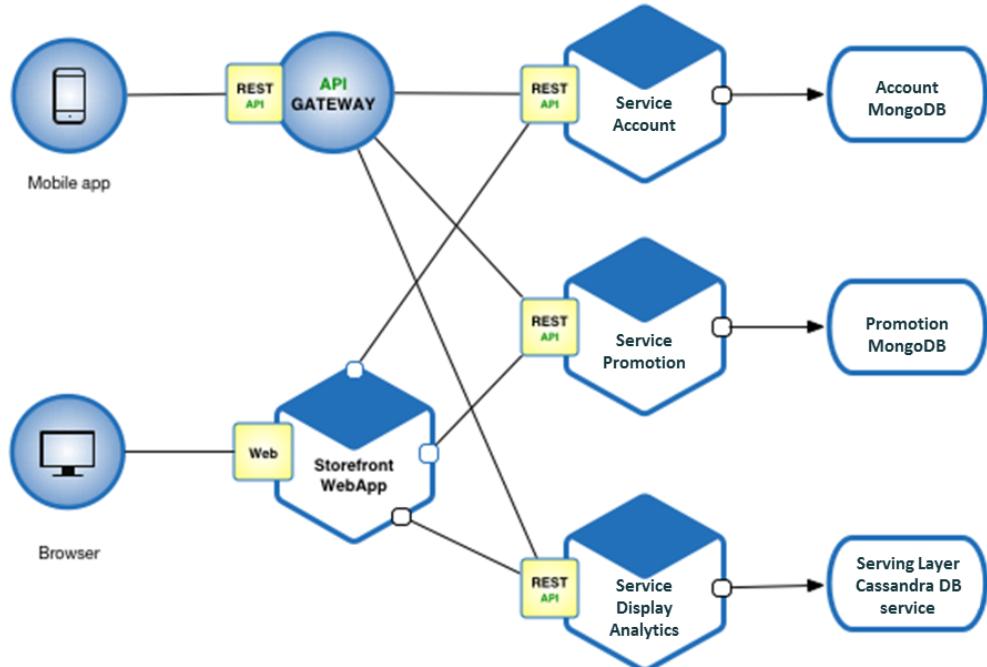


FIGURE 3.2. Sofrecom's subsidiaries and customers. Figure taken from [1].

3.1.1.3 Deployment architectural design

We have designed our deployment architecture showing all containers communication and linking pipelines. The next two figures show the big data and the microservices deployment architecture.

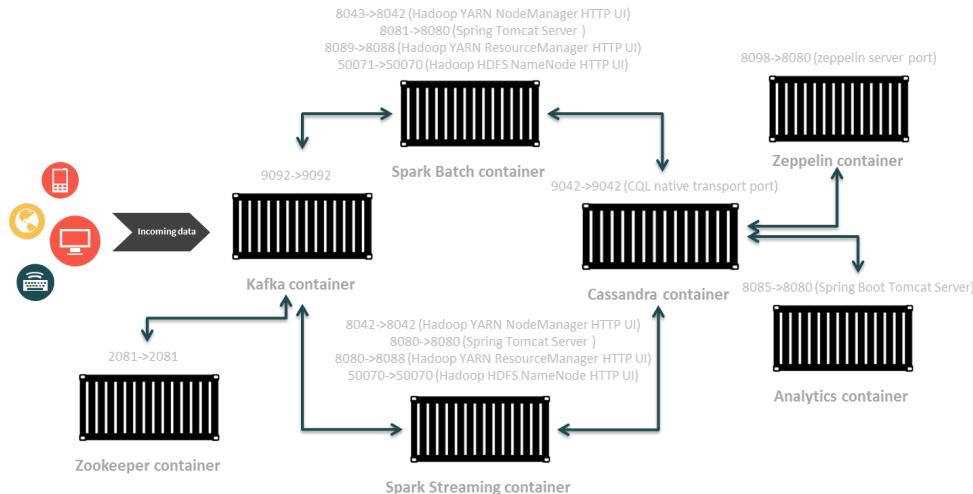


FIGURE 3.3. Sofrecom's subsidiaries and customers. Figure taken from [1].

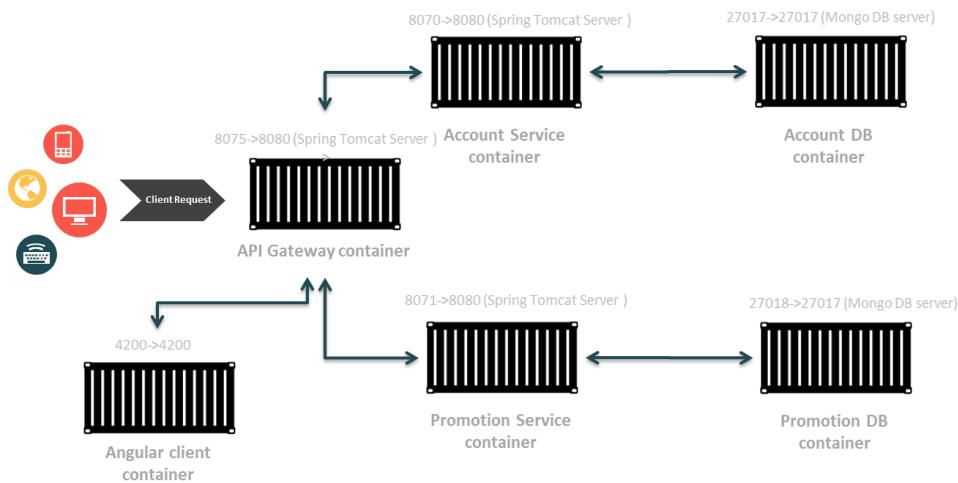


FIGURE 3.4. Sofrecom's subsidiaries and customers. Figure taken from [1].

3.1.1.4 Summary

Finally, the next figure resume all patterns and architecture described before.

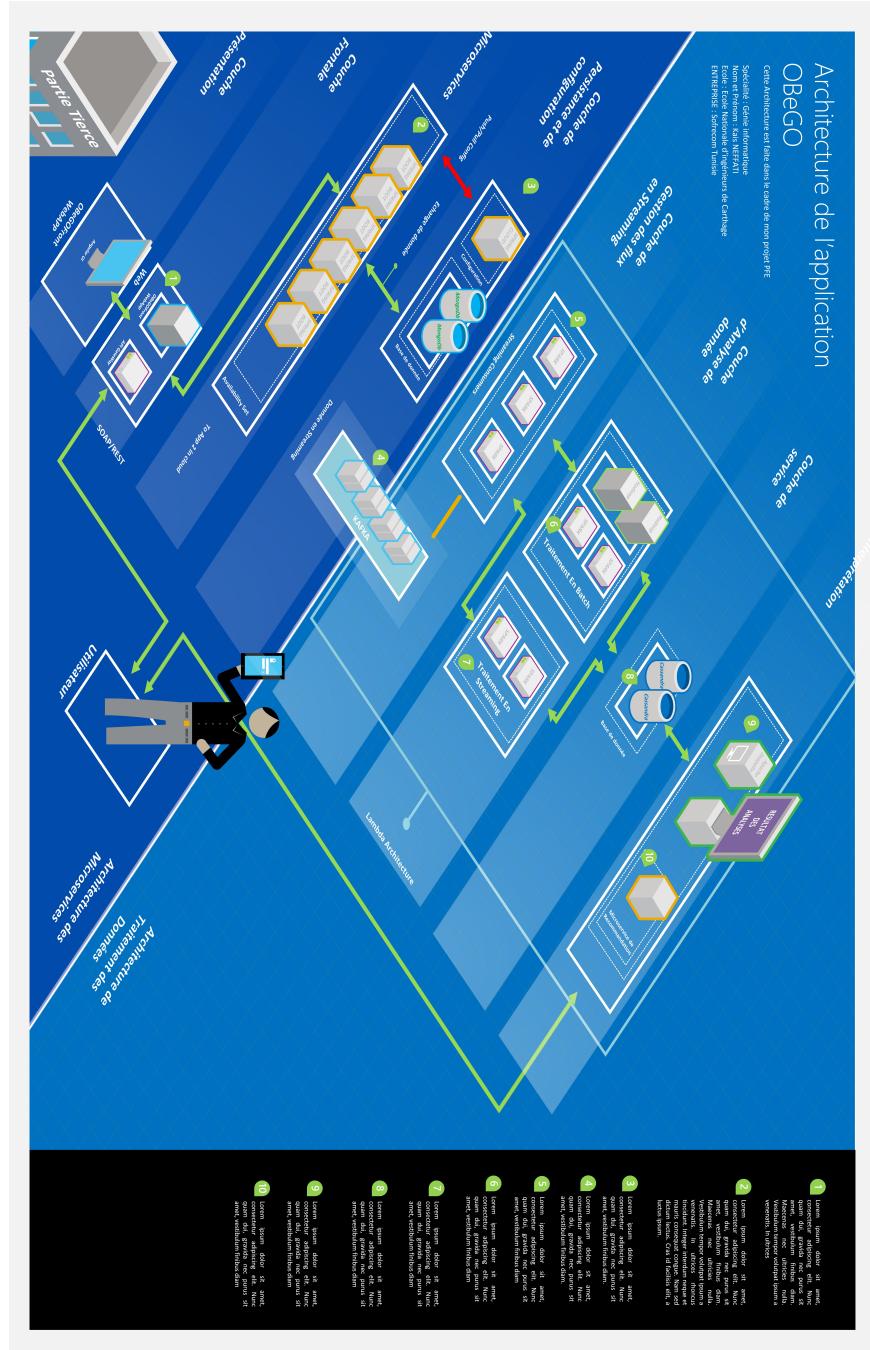


FIGURE 3.5. Sofrecom's subsidiaries and customers. Figure taken from [1].

3.2 Detailed design

This section will be devoted to the detailed design of our application, we start with the package diagram. Then we present the class diagram. And finally we model the dynamic aspect of the system using sequence diagrams

3.2.1 Package Diagram

A package diagram is a collection of static modeling elements (classes, packages, ...) that shows the structure of a model. It describes all the classes and their associations Packages in the

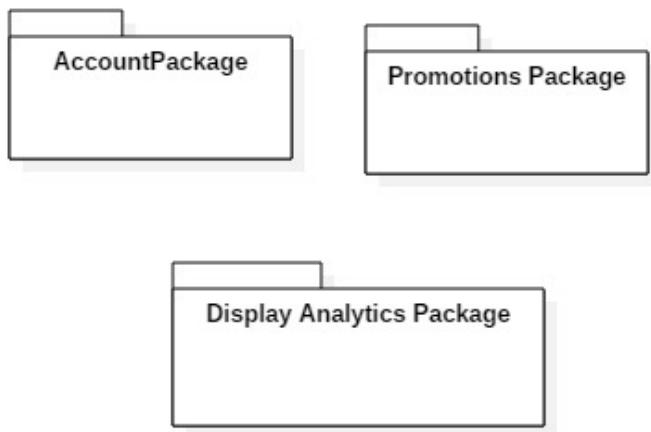


FIGURE 3.6. Sofrecom's subsidiaries and customers. Figure taken from [1].

microservices architecture are totally separated . This is one of the benifit of this approach. When creating a new fonctionaly just plug and pluck.

- Promotions package : this package handle promotions managment like creating , updating , searching or adding promotions.
- Account package : this package handle users managment like creating , updating , searching or adding promotions. This package expose a very important fonctionalty that allow updating user location. Users location are display in a heatMap.
- Display analytics package : This package takes care of data analytics results and display it in a dashboard. Allowing users to see realtime analytics results.

3.2.2 Class Diagram

After completing the preliminary design phase, we will now detail further this design by presenting the class diagrams of the different packages.

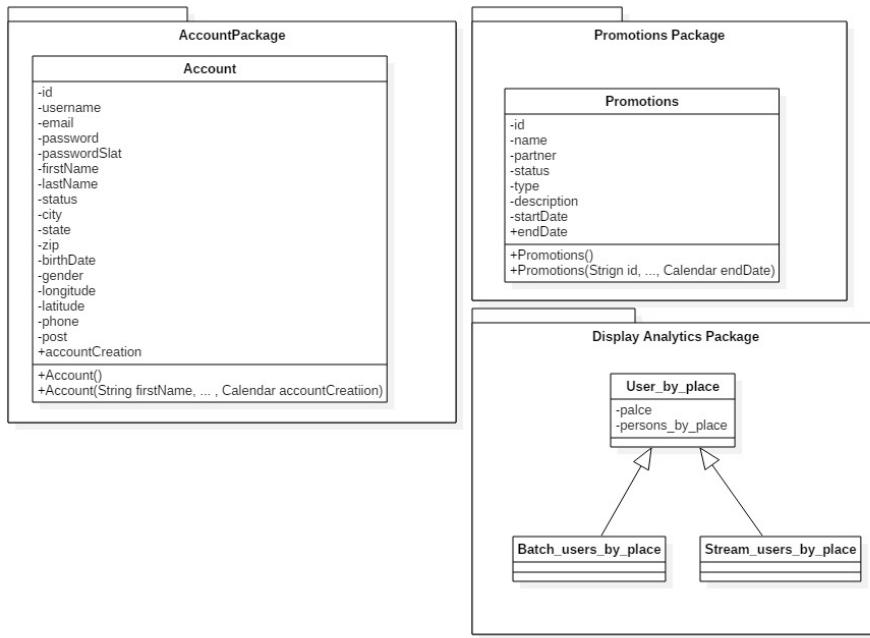


FIGURE 3.7. Sofrecom's subsidiaries and customers. Figure taken from [1].

- Account class : This class is used to define the user informations. It presents the Name, first name, email, login, passwords , essentially the longitude and latitude used for the heatMap chart ...
- Promotions class : This class is used to define the promotion informations. It presents the name, partner name, status, type, description ...
- Users by place class : this class is an abstract class that have two implementation which are Batch users by place and streaming users par place. It allows as to track users basing on locations and determine wich place is the most visited one. *For exmaple : if i am having a supermaket , by putting beacons in the each store corner , i will be mesuring the users vist by place in real time .*

3.2.3 Object Sequence Diagrams

A sequence diagram is an interaction diagram that shows how objects operate with one another and in what order. It is a construct of a message sequence chart. A sequence diagram shows object interactions arranged in time sequence.

3.2.3.1 Display analytics sequence diagrams

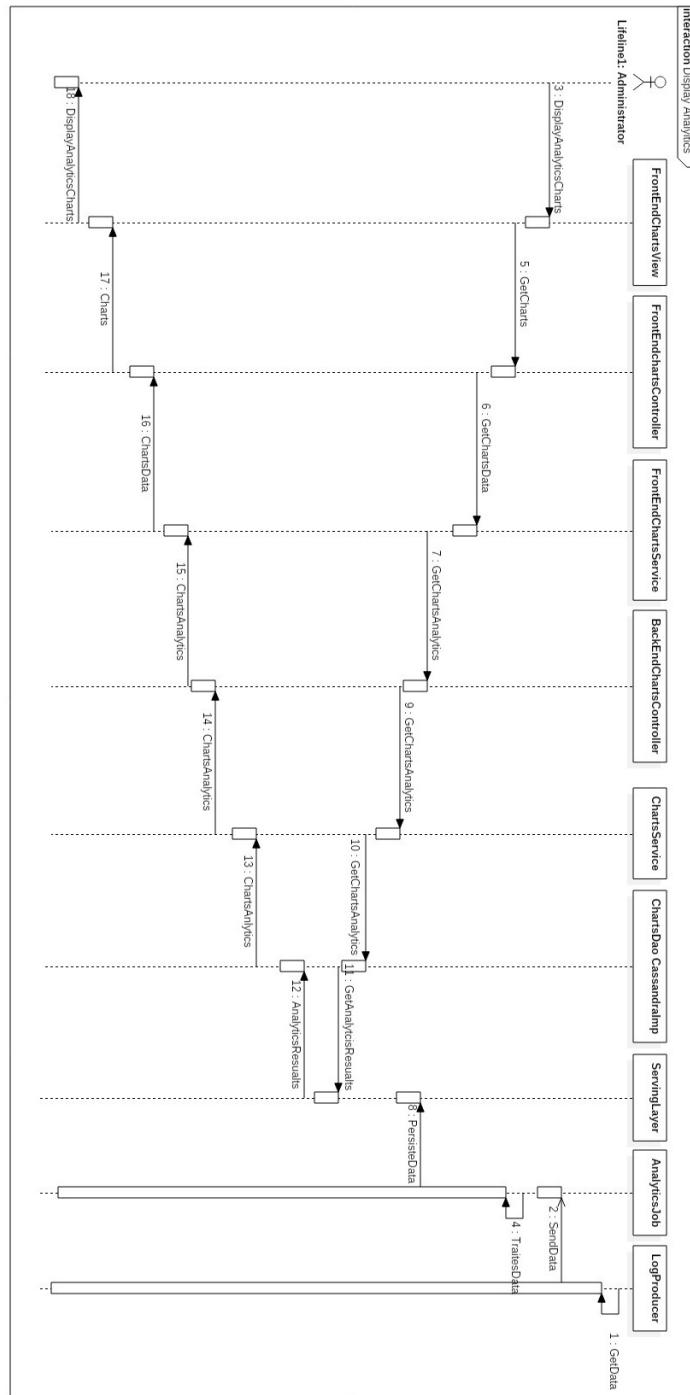


FIGURE 3.8. Sofrecom's subsidiaries and customers. Figure taken from [1].

3.2.3.2 Display analytics sequence diagrams

..

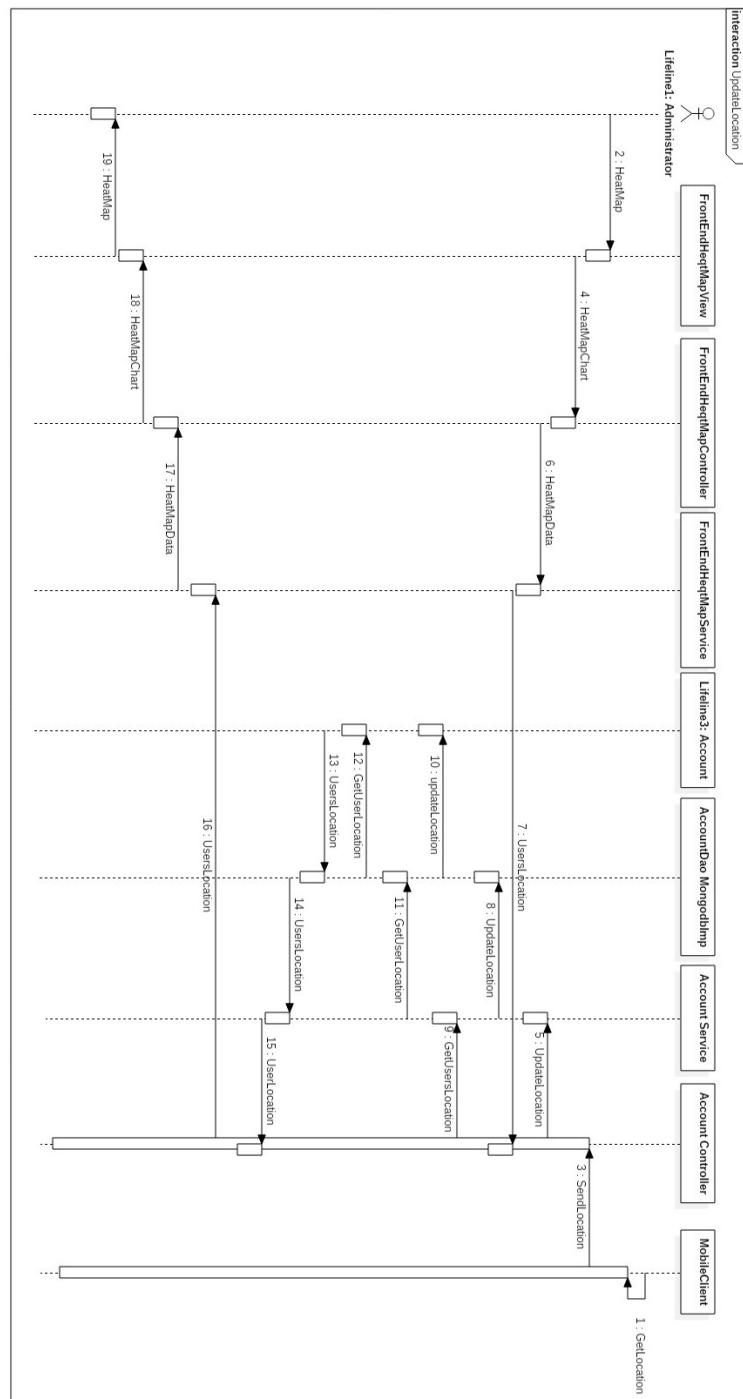


FIGURE 3.9. Sofrecom's subsidiaries and customers. Figure taken from [1].

3.3 Conclusion

Two main parts are waiting to be implemented, the first one is the lambda architecture and the data analytics strategy. the second one is the microservices approach and the application logic.

IMPLEMENTATION

In this section we will describe our project work. We will start by presenting the working environment . Then we move to talk about each part of our project , presenting screenshots and giving the examples of each Implementation.

4.1 Work environnement

This part presents the Hardware and Software environment available for this project.

4.1.1 Hardware environnement

In our project we developed our application using a computer with the following characteristics :

- Process Intel Core i5 4310M .
- Memory RAM: 16 GB .
- Space Disces: 500 GB.

4.2 Software environnement

We've used for our project the following Software :

- VMware Workstation : Our application is deployed under VMware Workstation virtualization platform as a set of decoupled service. Each service per container.
- Centos 7 : Our application operating system is Centos 7 carring Docker as demon.

- Docker : is our main virtualization tool and our containers manager.
- IntelliJ idea : is our application Back End IDE , we used it for Scala and java (spring) development.
- Visual Studio Code : is our application Front End IDE , we used it for angular 2 and webpack (TypeScript development).

4.3 Realized work

4.3.1 Environnement preparation

Our application require some containers preparation. The following list contain all needed containers.

- Kafka container : used for collecting data and stream them through a message queue. This container does not existe in dockerhub repository so we should Dockerize it(Dockerize kafka means create a container that contain kafka)
- Zookeeper container : used to save brokers address and messages offsets . this container exists in docker hub.
- Scala container : used to build the Hadoop container.This container does not existe in dockerhub repository so we should Dockerize it
- Hadoop container : used to build the spark container and provide HDFS persistance and YARN (ressource manager) for the spark jobs. This container does not existe with the Hadoop 2.7 and centos 7 version, so we should Dockerize it
- Spark container : used to to consume the data from kafka , process it then feed it to the serving layer (persisting results in cassandra).
- Cassandra container : used as a serving layer . this container exists in docker hub , we will be using the 2.2 cassandra version.
- Zeppelin container : used as a notebook to interpret Cassandra and spark then display the results in a dashboard.
- Spring boot containers : we need some containers that carry each one a microservice. This containers just need to have a JVM.
- Mongodb containers : each service is linked to a mongodb container. this containers exists in dockerhub.
- FrontEnd container : container for the angular dashboard.

4.3.1.1 Dockernize Tools

In this section we present the process of dockernizing our containers .All installation setps are written in Dockerfiles and each container has an OS kernel. We will be using alpine and centos 7. The next figure show all our docker images:

```

root@localhost obego]# docker images
REPOSITORY          TAG      IMAGE ID      CREATED       SIZE
neffatikais/analyticsmicroservice    latest   4669dd29b0da  9 weeks ago  212.8 MB
docker.io/mongo        latest   b7721338ba63  9 weeks ago  359.8 MB
docker.io/node        latest   329a17fdebdd  10 weeks ago 665.2 MB
neffatikais/zepelin     latest   47f3714af2d0  11 weeks ago 1.854 GB
neffatikais/spark       latest   c2aa2f61f086  11 weeks ago 1.714 GB
neffatikais/hadoop     latest   c5efd8dc6bb8  11 weeks ago 1.286 GB
neffatikais/scala       latest   b8e6fbff604c  11 weeks ago 933.3 MB
neffatikais/kafka       latest   8db617e2f00c  11 weeks ago 164.5 MB
docker.io/kong         latest   a9672f5ad94c  12 weeks ago 319.8 MB
docker.io/centos       7        98d35105a391  12 weeks ago 192.5 MB
docker.io/cassandra     2.2     2e99f7db0fd1  3 months ago 365.6 MB
docker.io/zookeeper     latest   3b83d9104a4c  3 months ago 128.5 MB
docker.io/frolvlad/alpine-oraclejdk8  slim    00d8610f052e  3 months ago 166.6 MB
docker.io/alpine        3.5     4a415e366388  3 months ago 3.984 MB
[ root@localhost obego]#

```

FIGURE 4.1. Docker images

Dockernize Kafka To dockernize this container, we will start by putting a base OS version which is alpine linux . The next step is the Jdk installation. After fixing all paths we move to install kafka.



FIGURE 4.2. Kafka container

Dockernize Scala To dockernize this container, we will start by putting a base OS version which is centos 7 . The next step is the Jdk installation. After fixing all paths we move to install Scala.



FIGURE 4.3. Scala container

Dockernize Hadoop To dockernize this container, we will start from the previous Scala container. The next step is the Hadoop installation which takes too much time due to the big configuration amount.



FIGURE 4.4. Hadoop Container

4.3. REALIZED WORK

After dockernizing Hadoop 2.7 we are able now to see some Dashboard. Hadoop YARN ResourceManager HTTP UI :

ID	User	Name	Application Type	Queue	StartTime	FinishTime	State	FinalStatus	Progress
application_1491994560265_0001	root	StreamingProcessing.Streaming	SPARK	default	Wed Apr 12 06:57:39 -0400 2017	N/A	ACCEPTED	UNDEFINED	

FIGURE 4.5. Hadoop YARN ResourceManager HTTP UI

Hadoop HDFS NameNode HTTP UI :

Started:	Tue May 09 13:43:40 UTC 2017
Version:	2.7.3. rbaa81f7ce8c952be5982de4719c1cb951ccff
Compiled:	2016-08-17T01:41Z by root from branch-2.7.3
Cluster ID:	CD-9d24b749-b266-433c-b3f4-7a7fc7fb655
Block Pool ID:	BP-1443018692-172.17.0.2-1490287481800

FIGURE 4.6. Hadoop HDFS NameNode HTTP UI

We can Check Files existance through the HDFS NameNode UI :

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
drwxrwxr-x	root	supergroup	0 B	3/30/2017, 5:59:56 AM	0	0 B	timestamp_hour=1490864400000
drwxrwxr-x	root	supergroup	0 B	3/30/2017, 5:59:56 AM	0	0 B	timestamp_hour=1490868000000
drwxrwxr-x	root	supergroup	0 B	3/30/2017, 5:59:56 AM	0	0 B	timestamp_hour=1490871600000
drwxrwxr-x	root	supergroup	0 B	3/30/2017, 5:59:47 AM	0	0 B	timestamp_hour=1490875200000
drwxrwxr-x	root	supergroup	0 B	3/31/2017, 11:49:56 AM	0	0 B	timestamp_hour=1490972400000
drwxrwxr-x	root	supergroup	0 B	3/31/2017, 11:50:01 AM	0	0 B	timestamp_hour=1490976000000
drwxrwxr-x	root	supergroup	0 B	3/31/2017, 11:49:50 AM	0	0 B	timestamp_hour=1490979600000
drwxrwxr-x	root	supergroup	0 B	3/31/2017, 11:49:40 AM	0	0 B	timestamp_hour=1490983200000

FIGURE 4.7. Browser files HDFS

Dockernize Spark This container is the main container. We built the last two ones for this step. This container is based on the haddop one . Spark uses HDFS to persist data and uses YARN as a ResourceManager. This container require much more then the last one in terms of configuration. We chose to built spark on yarn to benefit the real sparks power.

So spark present in totaly 3 deployment mode :

Standalone deployment: With the standalone deployment one we can statically allocate resources on all or a subset of machines in a Hadoop cluster and run Spark side by side with Hadoop MR. The user can then run arbitrary Spark jobs on his HDFS data. Its simplicity makes this the deployment of choice for many Hadoop 1.x users.

Hadoop Yarn deployment: Hadoop users who have already deployed or are planning to deploy Hadoop Yarn can simply run Spark on YARN without any pre- installation or administrative access required. This allows users to easily integrate Spark in their Hadoop stack and take advantage of the full power of Spark, as well as of other components running on top of Spark.

Spark In MapReduce (SIMR): For the Hadoop users that are not running YARN yet, another option, in addition to the standalone deployment, is to use SIMR to launch Spark jobs inside MapReduce. With SIMR, users can start experimenting with Spark and use its shell within a couple of minutes after downloading it! This tremendously lowers the barrier of deployment, and lets virtually everyone play with Spark.

Notice: Spark is an in-Memory data processing framework so when trying to configure it make sure you put the right RM consumption.



FIGURE 4.8. Spark container .

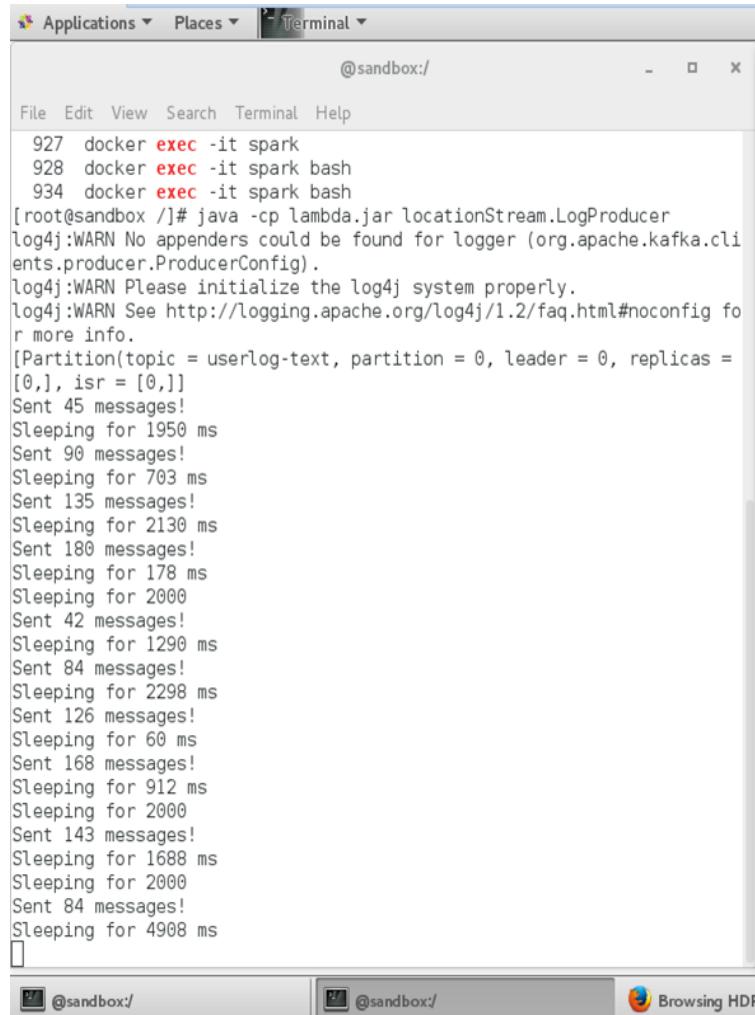
Dockernize Zeppelin Zeppelin is a notebook that allow as to interpret spark or cassandra . Zeppelin container is built on a an alpine linux operating system.



FIGURE 4.9. Zeppelin container.

4.3.2 Log Producer

The log producer is a Scala program that act as a producer for our kafka OSS. We ve used this log producer to perform a real data stream flow.



The screenshot shows a terminal window titled 'Terminal' with the command '@sandbox:/'. The window displays the output of a Scala application log producer. The logs show the application starting up, initializing log4j, and then sending 45 messages. It then enters a loop of sleeping for 1950 ms and sending 90 messages. This pattern repeats multiple times, with occasional variations in sleep duration and message count.

```
927 docker exec -it spark
928 docker exec -it spark bash
934 docker exec -it spark bash
[root@sandbox /]# java -cp lambda.jar locationStream.LogProducer
log4j:WARN No appenders could be found for logger (org.apache.kafka.clients.producer.ProducerConfig).
log4j:WARN Please initialize the log4j system properly.
log4j:WARN See http://logging.apache.org/log4j/1.2/faq.html#noconfig for more info.
[Partition(topic = userlog-text, partition = 0, leader = 0, replicas = [0,], isr = [0,])
Sent 45 messages!
Sleeping for 1950 ms
Sent 90 messages!
Sleeping for 703 ms
Sent 135 messages!
Sleeping for 2130 ms
Sent 180 messages!
Sleeping for 178 ms
Sleeping for 2000
Sent 42 messages!
Sleeping for 1290 ms
Sent 84 messages!
Sleeping for 2298 ms
Sent 126 messages!
Sleeping for 60 ms
Sent 168 messages!
Sleeping for 912 ms
Sleeping for 2000
Sent 143 messages!
Sleeping for 1688 ms
Sleeping for 2000
Sent 84 messages!
Sleeping for 4908 ms
```

FIGURE 4.10. Scala application log producer.

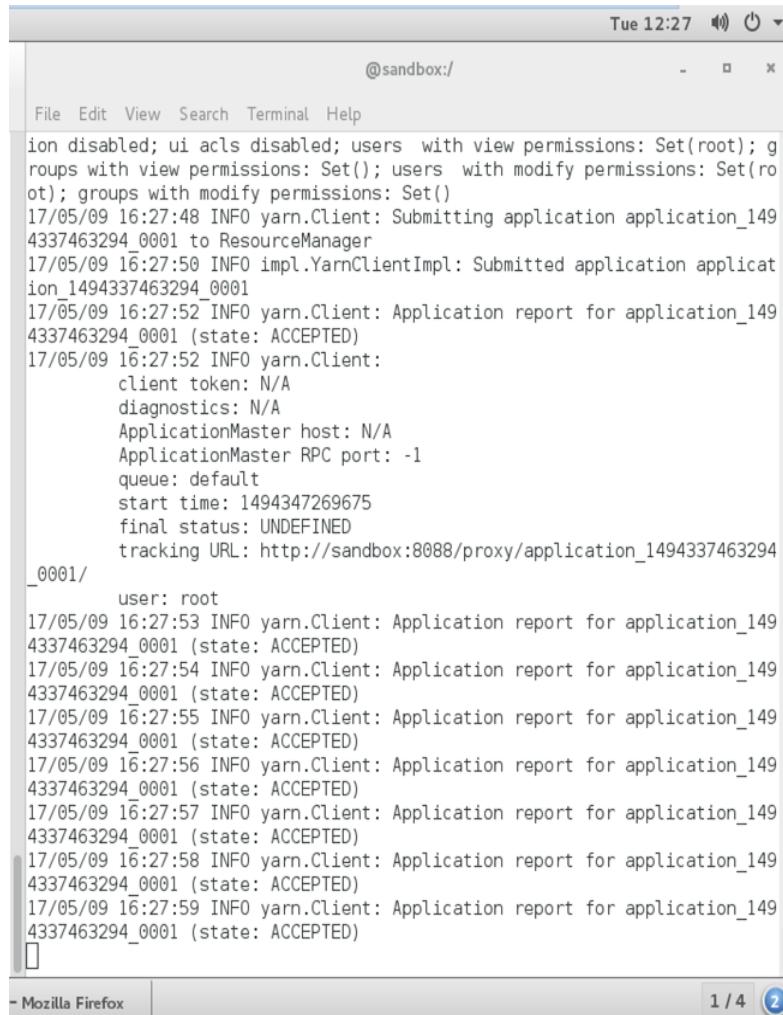
4.3.3 Spark Jobs

Applying the lambda architecture require running at least two jobs, in our case we've created three spark jobs using Scala :

- **BatchSaveTOHDFS:** This job persists realtime data to HDFS. It consumes data flow from Kafka and persist it to HDFS continuously.
- **StreamingProcessing:** This job is the realtime engine, it consumes data from Kafka and process it and gives results in realtime.
- **BatchProcessing:** This job is the batch engine, it processes the data already persisted in HDFS and gives results each 8 hours.

4.3.3.1 BatchProcessing

The next figure show the BatchProcessing job in it's accepted status:



The screenshot shows a terminal window titled 'Tue 12:27' with a progress bar icon. The window title is '@sandbox:/'. The terminal interface includes standard menu options: File, Edit, View, Search, Terminal, Help. Below the menu is a toolbar with icons for minimize, maximize, and close. The main area displays a log of YARN Client activity for application_1494337463294_0001. The log entries are as follows:

```
ion disabled; ui acls disabled; users with view permissions: Set(root); g
roups with view permissions: Set(); users with modify permissions: Set(ro
ot); groups with modify permissions: Set()
17/05/09 16:27:48 INFO yarn.Client: Submitting application application_149
4337463294_0001 to ResourceManager
17/05/09 16:27:50 INFO impl.YarnClientImpl: Submitted application applicat
ion_1494337463294_0001
17/05/09 16:27:52 INFO yarn.Client: Application report for application_149
4337463294_0001 (state: ACCEPTED)
17/05/09 16:27:52 INFO yarn.Client:
    client token: N/A
    diagnostics: N/A
    ApplicationMaster host: N/A
    ApplicationMaster RPC port: -1
    queue: default
    start time: 1494347269675
    final status: UNDEFINED
    tracking URL: http://sandbox:8088/proxy/application_1494337463294
_0001/
    user: root
17/05/09 16:27:53 INFO yarn.Client: Application report for application_149
4337463294_0001 (state: ACCEPTED)
17/05/09 16:27:54 INFO yarn.Client: Application report for application_149
4337463294_0001 (state: ACCEPTED)
17/05/09 16:27:55 INFO yarn.Client: Application report for application_149
4337463294_0001 (state: ACCEPTED)
17/05/09 16:27:56 INFO yarn.Client: Application report for application_149
4337463294_0001 (state: ACCEPTED)
17/05/09 16:27:57 INFO yarn.Client: Application report for application_149
4337463294_0001 (state: ACCEPTED)
17/05/09 16:27:58 INFO yarn.Client: Application report for application_149
4337463294_0001 (state: ACCEPTED)
17/05/09 16:27:59 INFO yarn.Client: Application report for application_149
4337463294_0001 (state: ACCEPTED)
```

The bottom of the terminal window shows the Mozilla Firefox logo and the page number '1 / 4'.

FIGURE 4.11. BatchProcessing spark job in accepted status.

CHAPTER 4. IMPLEMENTATION

The next figure show the BatchProcessing job in it's running status:

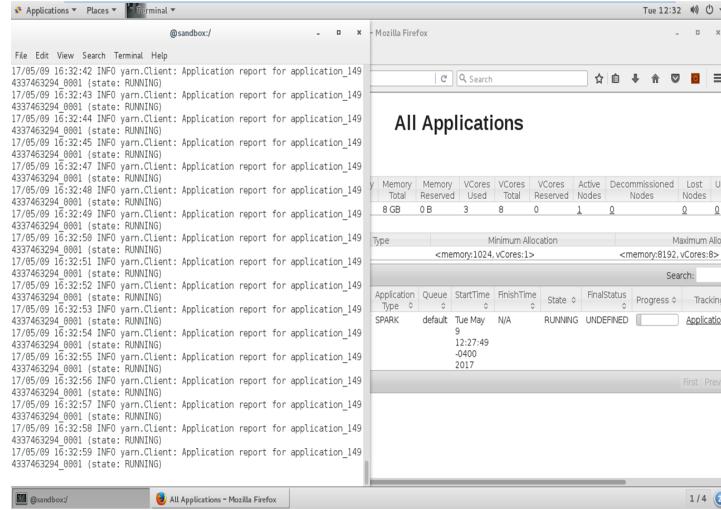


FIGURE 4.12. BatchProcessing spark job in running status.

4.3.3.2 StreamingProcessing

The next figure show the StreamingProcessing job in it's accepted status:

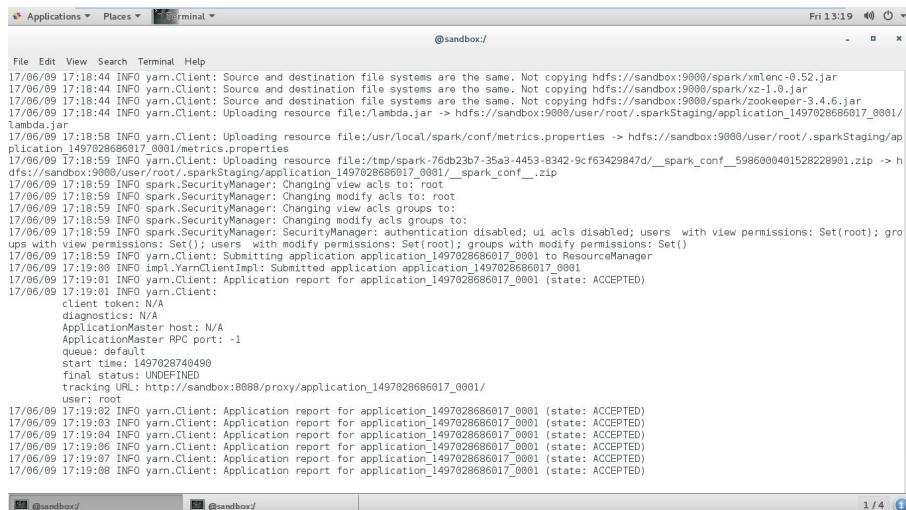
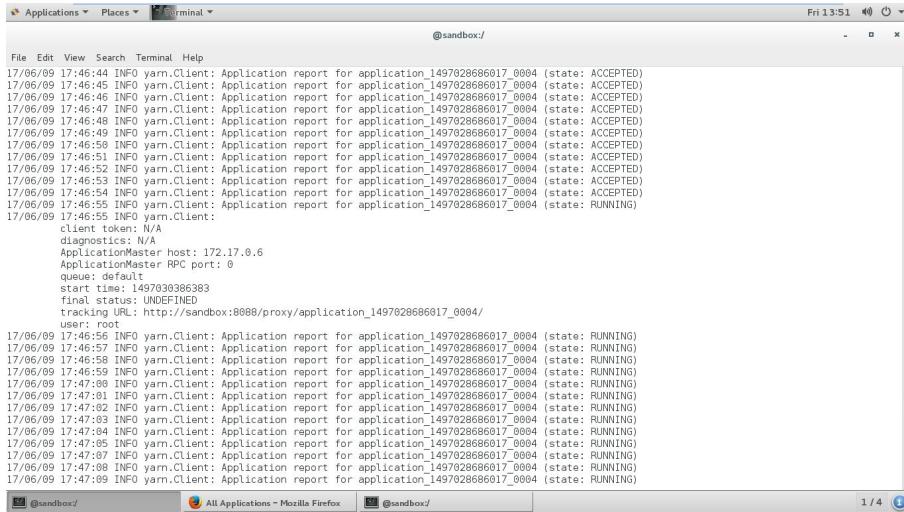


FIGURE 4.13. StreamingProcessing spark job in accepted status.

4.3. REALIZED WORK

The next figure show the BatchProcessing job in it's running status:



```

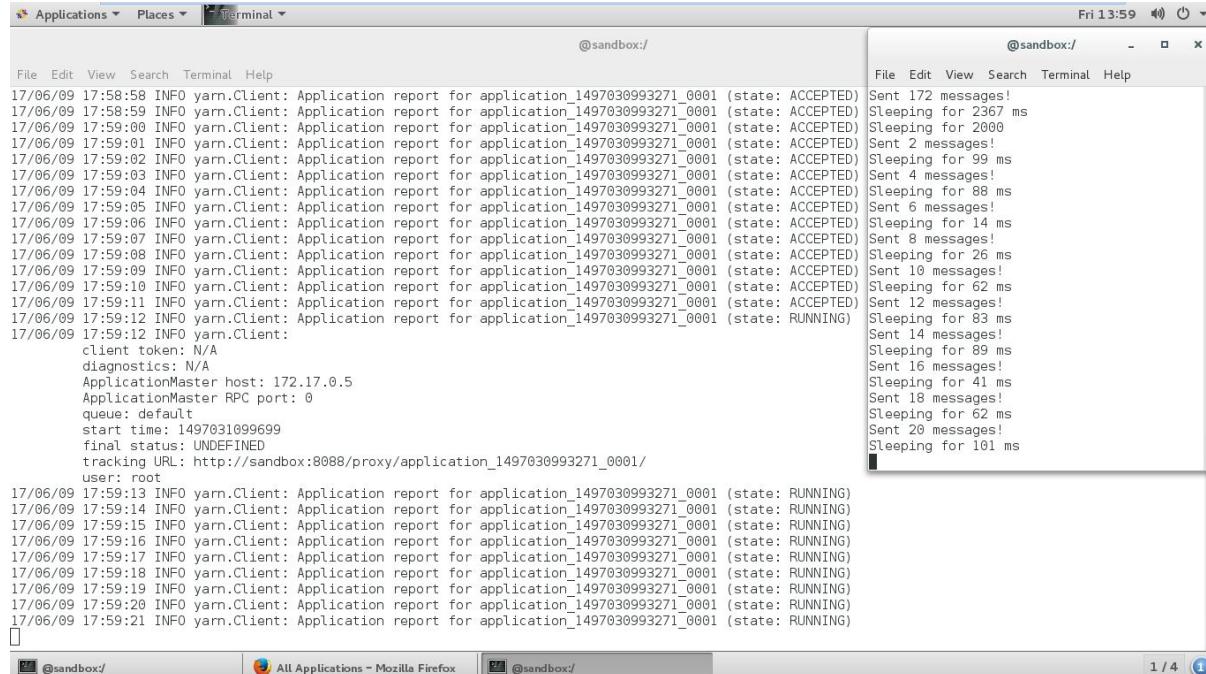
File Edit View Search Terminal Help
@sandbox:/ Fri 13:51
17/06/09 17:46:44 INFO yarn.Client: Application report for application_1497028686017_0004 (state: ACCEPTED)
17/06/09 17:46:45 INFO yarn.Client: Application report for application_1497028686017_0004 (state: ACCEPTED)
17/06/09 17:46:46 INFO yarn.Client: Application report for application_1497028686017_0004 (state: ACCEPTED)
17/06/09 17:46:47 INFO yarn.Client: Application report for application_1497028686017_0004 (state: ACCEPTED)
17/06/09 17:46:48 INFO yarn.Client: Application report for application_1497028686017_0004 (state: ACCEPTED)
17/06/09 17:46:49 INFO yarn.Client: Application report for application_1497028686017_0004 (state: ACCEPTED)
17/06/09 17:46:50 INFO yarn.Client: Application report for application_1497028686017_0004 (state: ACCEPTED)
17/06/09 17:46:51 INFO yarn.Client: Application report for application_1497028686017_0004 (state: ACCEPTED)
17/06/09 17:46:52 INFO yarn.Client: Application report for application_1497028686017_0004 (state: ACCEPTED)
17/06/09 17:46:53 INFO yarn.Client: Application report for application_1497028686017_0004 (state: ACCEPTED)
17/06/09 17:46:54 INFO yarn.Client: Application report for application_1497028686017_0004 (state: ACCEPTED)
17/06/09 17:46:55 INFO yarn.Client: Application report for application_1497028686017_0004 (state: RUNNING)
17/06/09 17:46:55 INFO yarn.Client: Client token: N/A
diagnostics: N/A
ApplicationMaster host: 172.17.0.6
ApplicationMaster RPC port: 0
queue: default
start time: 1497030986383
final status: UNDEFINED
tracking URL: http://sandbox:8088/proxy/application_1497028686017_0004/
user: root
17/06/09 17:46:56 INFO yarn.Client: Application report for application_1497028686017_0004 (state: RUNNING)
17/06/09 17:46:57 INFO yarn.Client: Application report for application_1497028686017_0004 (state: RUNNING)
17/06/09 17:46:58 INFO yarn.Client: Application report for application_1497028686017_0004 (state: RUNNING)
17/06/09 17:46:59 INFO yarn.Client: Application report for application_1497028686017_0004 (state: RUNNING)
17/06/09 17:47:00 INFO yarn.Client: Application report for application_1497028686017_0004 (state: RUNNING)
17/06/09 17:47:01 INFO yarn.Client: Application report for application_1497028686017_0004 (state: RUNNING)
17/06/09 17:47:02 INFO yarn.Client: Application report for application_1497028686017_0004 (state: RUNNING)
17/06/09 17:47:03 INFO yarn.Client: Application report for application_1497028686017_0004 (state: RUNNING)
17/06/09 17:47:04 INFO yarn.Client: Application report for application_1497028686017_0004 (state: RUNNING)
17/06/09 17:47:05 INFO yarn.Client: Application report for application_1497028686017_0004 (state: RUNNING)
17/06/09 17:47:07 INFO yarn.Client: Application report for application_1497028686017_0004 (state: RUNNING)
17/06/09 17:47:08 INFO yarn.Client: Application report for application_1497028686017_0004 (state: RUNNING)
17/06/09 17:47:09 INFO yarn.Client: Application report for application_1497028686017_0004 (state: RUNNING)

```

FIGURE 4.14. StreamingProcessing spark job in running status.

4.3.3.3 BatchSaveTOHDFS

The next figure show the BatchSaveToHDFS job in it's accepted and running status:



```

File Edit View Search Terminal Help
@sandbox:/ Fri 13:59
17/06/09 17:58:58 INFO yarn.Client: Application report for application_1497030993271_0001 (state: ACCEPTED)
17/06/09 17:58:59 INFO yarn.Client: Application report for application_1497030993271_0001 (state: ACCEPTED)
17/06/09 17:59:00 INFO yarn.Client: Application report for application_1497030993271_0001 (state: ACCEPTED)
17/06/09 17:59:01 INFO yarn.Client: Application report for application_1497030993271_0001 (state: ACCEPTED)
17/06/09 17:59:02 INFO yarn.Client: Application report for application_1497030993271_0001 (state: ACCEPTED)
17/06/09 17:59:03 INFO yarn.Client: Application report for application_1497030993271_0001 (state: ACCEPTED)
17/06/09 17:59:04 INFO yarn.Client: Application report for application_1497030993271_0001 (state: ACCEPTED)
17/06/09 17:59:05 INFO yarn.Client: Application report for application_1497030993271_0001 (state: ACCEPTED)
17/06/09 17:59:06 INFO yarn.Client: Application report for application_1497030993271_0001 (state: ACCEPTED)
17/06/09 17:59:07 INFO yarn.Client: Application report for application_1497030993271_0001 (state: ACCEPTED)
17/06/09 17:59:08 INFO yarn.Client: Application report for application_1497030993271_0001 (state: ACCEPTED)
17/06/09 17:59:09 INFO yarn.Client: Application report for application_1497030993271_0001 (state: ACCEPTED)
17/06/09 17:59:10 INFO yarn.Client: Application report for application_1497030993271_0001 (state: ACCEPTED)
17/06/09 17:59:11 INFO yarn.Client: Application report for application_1497030993271_0001 (state: ACCEPTED)
17/06/09 17:59:12 INFO yarn.Client: Application report for application_1497030993271_0001 (state: RUNNING)
17/06/09 17:59:12 INFO yarn.Client: client token: N/A
diagnostics: N/A
ApplicationMaster host: 172.17.0.5
ApplicationMaster RPC port: 0
queue: default
start time: 1497031099699
final status: UNDEFINED
tracking URL: http://sandbox:8088/proxy/application_1497030993271_0001/
user: root
17/06/09 17:59:13 INFO yarn.Client: Application report for application_1497030993271_0001 (state: RUNNING)
17/06/09 17:59:14 INFO yarn.Client: Application report for application_1497030993271_0001 (state: RUNNING)
17/06/09 17:59:15 INFO yarn.Client: Application report for application_1497030993271_0001 (state: RUNNING)
17/06/09 17:59:16 INFO yarn.Client: Application report for application_1497030993271_0001 (state: RUNNING)
17/06/09 17:59:17 INFO yarn.Client: Application report for application_1497030993271_0001 (state: RUNNING)
17/06/09 17:59:18 INFO yarn.Client: Application report for application_1497030993271_0001 (state: RUNNING)
17/06/09 17:59:19 INFO yarn.Client: Application report for application_1497030993271_0001 (state: RUNNING)
17/06/09 17:59:20 INFO yarn.Client: Application report for application_1497030993271_0001 (state: RUNNING)
17/06/09 17:59:21 INFO yarn.Client: Application report for application_1497030993271_0001 (state: RUNNING)

File Edit View Search Terminal Help
@sandbox:/ Fri 13:59
Sent 172 messages!
Sleeping for 2367 ms
Sleeping for 2000
Sent 2 messages!
Sleeping for 99 ms
Sent 4 messages!
Sleeping for 88 ms
Sent 6 messages!
Sleeping for 14 ms
Sent 8 messages!
Sleeping for 26 ms
Sent 10 messages!
Sleeping for 62 ms
Sent 12 messages!
Sleeping for 83 ms
Sent 14 messages!
Sleeping for 89 ms
Sent 16 messages!
Sleeping for 41 ms
Sent 18 messages!
Sleeping for 62 ms
Sent 20 messages!
Sleeping for 101 ms

```

FIGURE 4.15. BatchSaveToHDFS spark job in accepted status.

CHAPTER 4. IMPLEMENTATION

The next figure show the job Hadoop Yarn UI status:

The screenshot shows the 'All Applications' page of the Hadoop Yarn UI. The title bar says 'All Applications - Mozilla Firefox'. The main content area is titled 'All Applications' with a 'hadoop' logo. On the left, there's a sidebar with a 'Cluster' section containing links like 'About', 'Nodes', 'Node Labels', 'Applications' (which is expanded), and 'Scheduler'. Below these are status counts: NEW (0), NEW_SAVING (0), SUBMITTED (0), ACCEPTED (0), RUNNING (1), FINISHED (0), FAILED (0), KILLED (0), and Scheduler. There's also a 'Tools' section. The main table is titled 'Cluster Metrics' and shows columns for Apps Submitted, Apps Pending, Apps Running, Apps Completed, Containers, Memory Running, Memory Used, Memory Total, Memory Reserved, Vcores Used, Vcores Total, Vcores Reserved, Active Nodes, and Decommissioned Nodes. The table has one row with values: Apps Submitted: 1, Apps Pending: 0, Apps Running: 1, Apps Completed: 0, Containers: 3, Memory Running: 6 GB, Memory Used: 8 GB, Memory Total: 0 B, Memory Reserved: 0 B, Vcores Used: 3, Vcores Total: 8, Vcores Reserved: 0, Active Nodes: 1, Decommissioned Nodes: 0. Below this is a 'Scheduler Metrics' table with columns for Scheduler Type, Scheduling Resource Type, and Minimum Allocation. It shows Capacity Scheduler, [MEMORY], <memory:1024, vCores:1>, and <memory:8192, vCores:1>. A table below lists applications with columns ID, User, Name, Application Type, Queue, Start Time, Finish Time, State, and Final Status. One entry is shown: application_1497030993271_0001, root, StreamingProcessing.StreamingSaveToHDFS, SPARK, default, Fri Jun 9 13:58:19 -0400 2017, RUNNING, UNDEFINED. At the bottom, it says 'Showing 1 to 1 of 1 entries'.

FIGURE 4.16. Running job status Hadoop Yarn UI .

The next figure show the Running application Master:

The screenshot shows the 'Application application_1497030993271_0001' dashboard. The title bar says 'Application application_1497030993271_0001 - Mozilla Firefox'. The main content area is titled 'Application application_1497030993271_0001' with a 'hadoop' logo. On the left, there's a sidebar with a 'Cluster' section containing links like 'About', 'Nodes', 'Node Labels', 'Applications' (which is expanded), and 'Scheduler'. Below these are status counts: NEW (0), NEW_SAVING (0), SUBMITTED (0), ACCEPTED (0), RUNNING (1), FINISHED (0), FAILED (0), KILLED (0), and Scheduler. There's also a 'Tools' section. The main area has two tabs: 'Kill Application' (selected) and 'Application Overview'. The 'Kill Application' tab shows details for the application: User: root, Name: StreamingProcessing.StreamingSaveToHDFS, Application Type: SPARK, Application Tags: YarnApplicationState: RUNNING: AM has registered with RM and started running., Queue: default, FinalStatus Reported by AM: Application has not completed yet., Started: Fri Jun 09 17:58:19 +0000 2017, Elapsed: 2mins, 38sec, Tracking URL: ApplicationMaster, and Diagnostics: . The 'Application Overview' tab shows Application Metrics: Total Resource Preempted: <memory:0, vCores:0>, Total Number of Non-AM Containers preempted: 0, Total Number of AM Containers Preempted: 0, Resource Preempted from Current Attempt: <memory:0, vCores:0>, Number of Non-AM Containers Preempted from Current Attempt: 0, and Aggregate Resource Allocation: 760627 MB-seconds, 370 vcore-seconds. At the bottom, it says 'Showing 1 to 1 of 1 matches'.

FIGURE 4.17. Application Master dashboard .

4.3. REALIZED WORK

The next figure show the running job log:

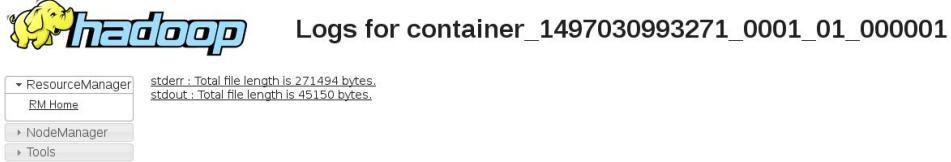


FIGURE 4.18. Job log .

The next figure show the running job output:

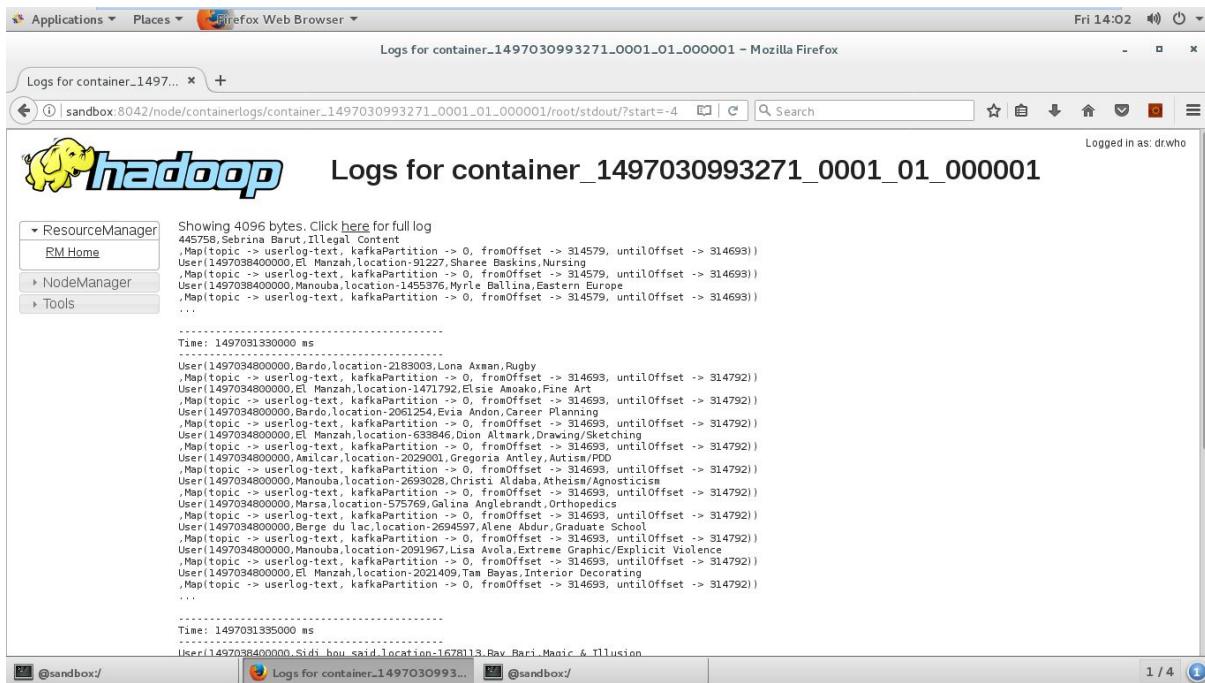


FIGURE 4.19. Job log Results .

CHAPTER 4. IMPLEMENTATION

4.3.4 Show Results on Apache Zeppelin

The next figure show Analytics results displayed on Zeppelin dashboard :

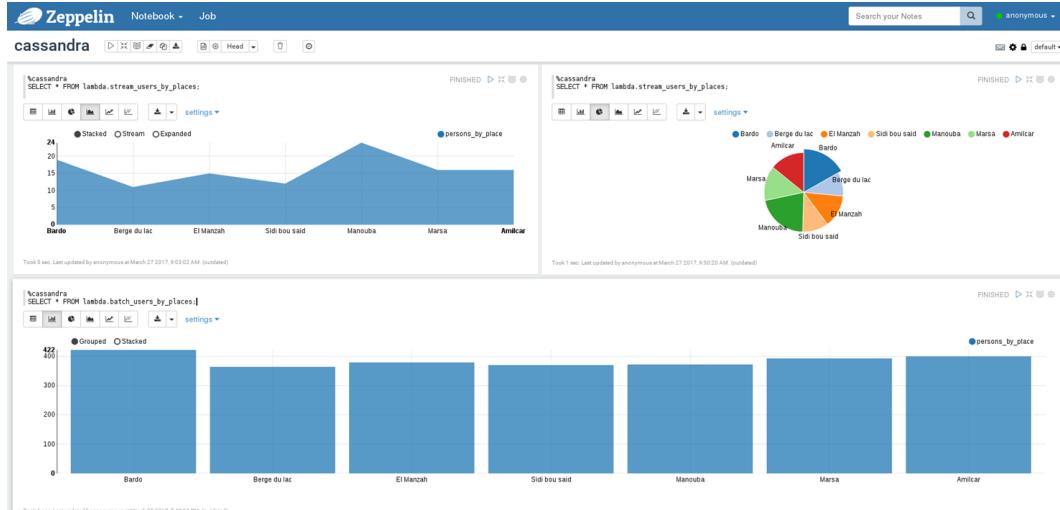


FIGURE 4.20. Zeppelin dashboard.

4.3.5 Account Management

4.3.5.1 Accounts

The next figure show the list Accounts :

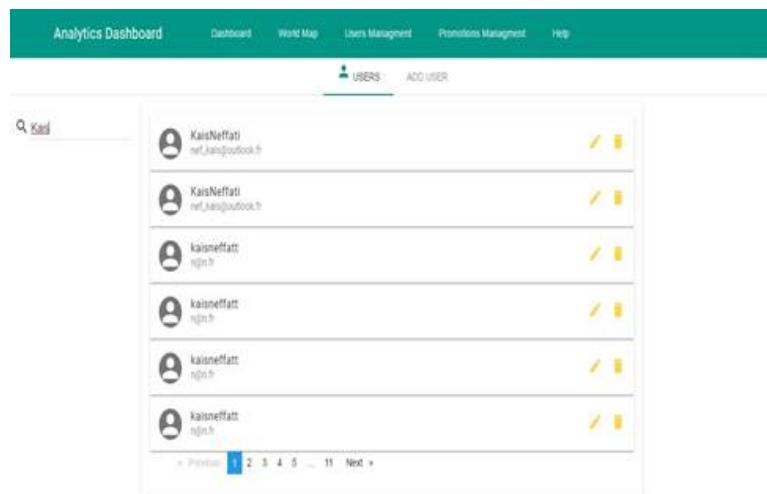
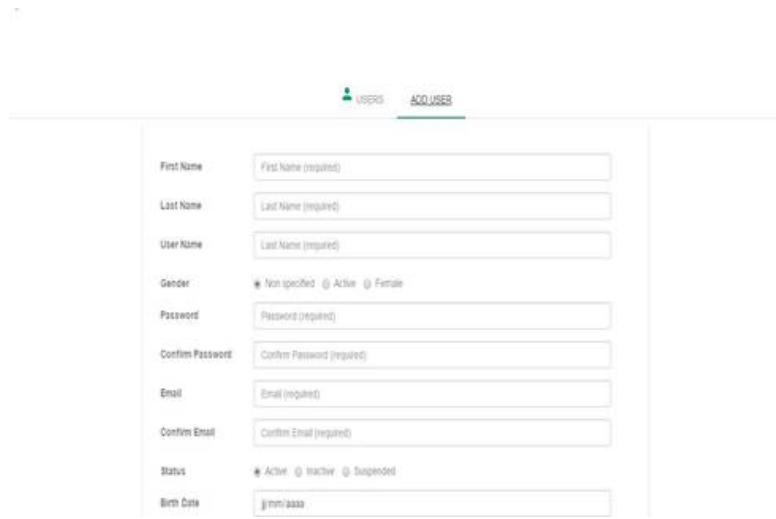


FIGURE 4.21. Account list.

4.3.5.2 Account add

The next figure show the form add account :



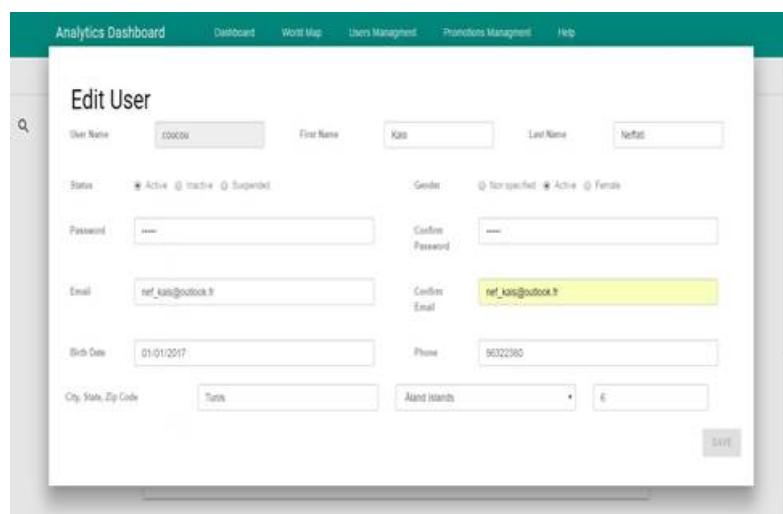
The screenshot shows a 'User Management' interface with a 'ADD USER' tab selected. The form contains the following fields:

First Name	First Name (required)
Last Name	Last Name (required)
User Name	User Name (required)
Gender	<input checked="" type="radio"/> Non specified <input type="radio"/> Active <input type="radio"/> Female
Password	Password (required)
Confirm Password	Confirm Password (required)
Email	Email (required)
Confirm Email	Confirm Email (required)
Status	<input checked="" type="radio"/> Active <input type="radio"/> Inactive <input type="radio"/> Suspended
Birth Date	j/mm/aaaa

FIGURE 4.22. Add account.

4.3.5.3 Account Update

The next figure show the form update account :



The screenshot shows an 'Edit User' dialog box. The user name is set to 'EDOUARD'. The form fields include:

User Name	EDOUARD	First Name	Kao	Last Name	Natali
Status	<input checked="" type="radio"/> Active <input type="radio"/> Inactive <input type="radio"/> Suspended	Gender	<input type="radio"/> Non specified <input checked="" type="radio"/> Active <input type="radio"/> Female		
Password	*****	Confirm Password	*****		
Email	ref_kao@outlook.fr	Confirm Email	ref_kao@outlook.fr		
Birth Date	01/01/2017	Phone	60322380		
City, State, Zip Code	Tours	Aland Islands	E		

At the bottom right of the dialog is a 'Save' button.

FIGURE 4.23. Update account.

4.3.6 Promotion Management

4.3.6.1 Promotions

The next figure show the list Promotions :

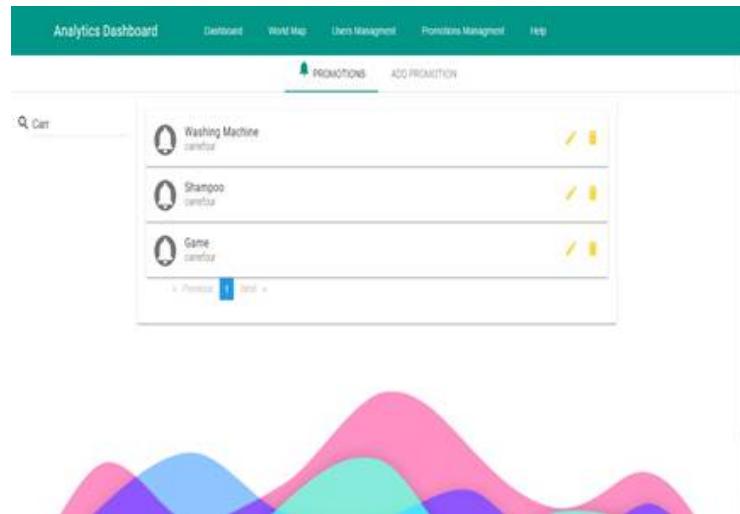


FIGURE 4.24. Promotions list.

4.3.6.2 Promotion add

The next figure show the form add Promotion :

The screenshot shows the 'ADD PROMOTION' form. The form fields are as follows: 'Name' (input field), 'Partner' (input field), 'Type' (input field), 'Status' (radio buttons for Active, Inactive, Suspended, with Active selected), 'Start Date' (input field), 'End Date' (input field), and 'description' (input field). A 'Save' button is located at the bottom right of the form.

FIGURE 4.25. Add promotions.

4.3. REALIZED WORK

4.3.6.3 Promotion Update

The next figure show the form update Promotion :

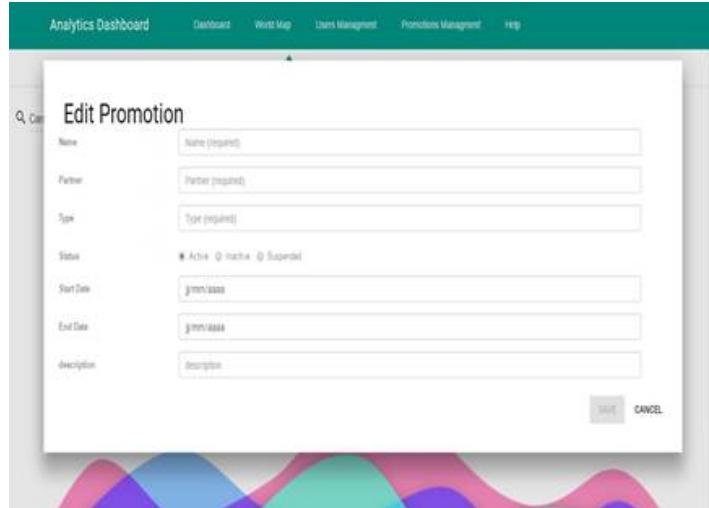


FIGURE 4.26. Update promotions.

4.3.7 Display Analytics

The next figures show the Data Analytics Doshboard , this dashboard is update each time new data comes :

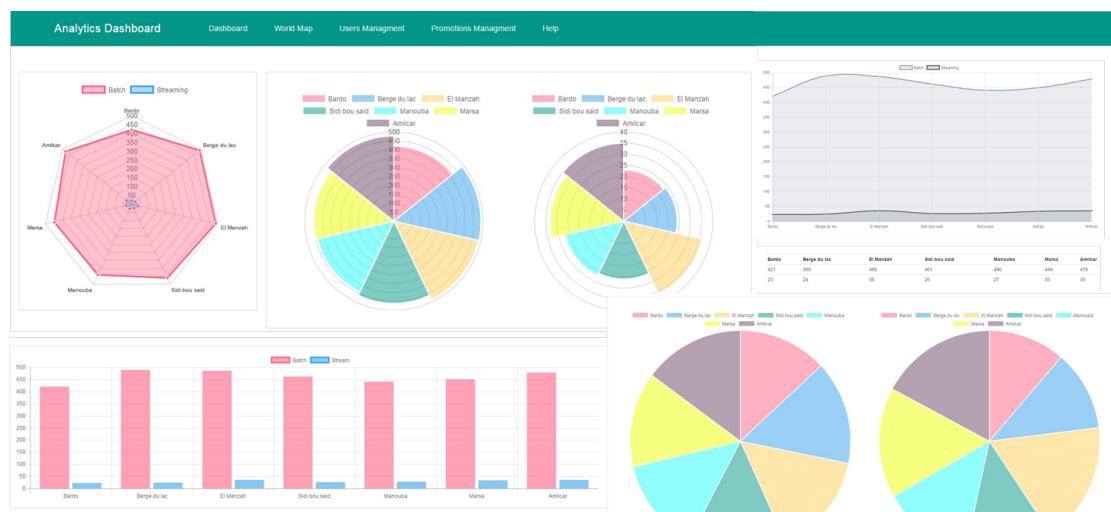


FIGURE 4.27. Display analytics dashboard.

4.3.8 Heat Map users

The next figure show the Data Analytics Doshboard , this dashboard is update each time new data comes :

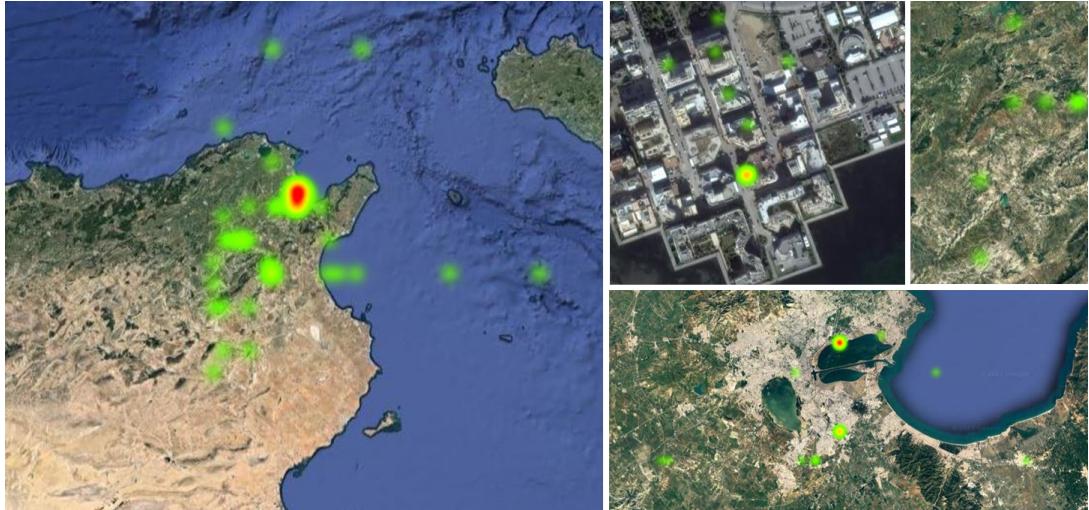


FIGURE 4.28. Heat map users dashboard.

4.3.9 Summary

In this chapter, we presented our project implementation. First, we have prepared the work environment which takes too much time configuring, adapting versions and dealing with system engineering issues. Second, we've moved to link each container to the other one. Then we started making some jobs like the log producer, the BatchSaveToHDFS, the StreamingProcessing and the Batch Processing. Third, we've implemented the microservices architecture that allows account (users) or promotions management. Finally, we have our results (Analytics, Cruds and Heat Map) displayed on a single dashboard.

4.4 General conclusion

This project was designed to handle scaling in term of features and treatment. Our project is still incomplete. There is a lot of features that we would like to involve. Some of them are:

- A gateway api, that handle authentication using Spring boot and zuul.
- A container manager like docker Swarm or Kubernities that takes care of high availability and scalability.
- So much functional features like recommendation engine using Spark mllib.
- A mobile application that takes the log producer place.
- ...

Finally, We have been focused on creating a well designed system that can handle all of the requirements presented in the second chapter. The technical and the functional branches are now implemented and move to gather in a straight line.

A P P E N D I X



APPENDIX A

Begins an appendix

BIBLIOGRAPHY

- [1] S. EQUIPE, *Sofrecom the know-how network*, (2012).

