

Master Project Report

August 30, 2016

GLAPP

Fabio Isler (09-115-965)
Man Wai Li (14-705-156)
Dinesh Pothineni (14-707-988)
Riccardo Patane (XX-XXX-XXX)

supervised by

Prof. Dr. Harald C. Gall

Dr. Philipp Leitner



**University of
Zurich** ^{UZH}



Abstract

Abstract

Introduction

Cloud computing has created a paradigm shift in the last few years, by making infrastructure available at lower costs and with higher efficiency of operations. These solutions are increasingly being adopted by enterprises and developers, as they can provision huge amount of resources to scale on-demand in order to meet their business needs. In cloud computing, resources such as CPU processing time, disk space, or networking capabilities, are rented and released as a service. Today, the most important model for delivering on the cloud promise is the Infrastructure-as-a-Service (IaaS) paradigm. In IaaS, virtual computing resources are acquired and released on demand, either via an Application Programming Interface (API) or web interface. Along with great flexibility of being able to get new resources on demand and pay for what you use, new problems arise. Selecting a cloud service provider can often be a quite challenging decision for the developer or a company, so is being able to monitor and evaluate these infrastructure resources on a regular basis. With so much of variance in cost and performance, it is imperative that one would look for a reliable deployment, monitoring solutions to strike a balance with application requirements. Furthermore, the complexity and skill required to manage multiple layers of application, data, middleware and operating systems can .

Understanding run time performance, behavior of various application components in real-time can enable us to take advantage of arbitrage opportunities that exist between different machines/regions or even other cloud providers. For example, an application can take advantage by moving closer to its users based on timezones or traffic to improve response time. Currently there is no flexibility to move freely between various cloud providers without great development effort and cost, however such an ability to move freely between providers enables us to benefit from cost and performance differences. We propose a cloud middleware that can not only take care of application deployment to the cloud, but also constantly monitor and trigger such necessary adaptations to benefit from these opportunities. This middleware also enables developers to specify their intended goals in terms of high level policies to govern the application behavior. The middleware can then break down these policies into low level objectives, in order to trigger adaptations by changing the state of application when required.

1.1 Global Living Cloud Applications

The aim of this project is to develop a platform for what we call *global living cloud applications* (GLAs). In a nutshell, GLAs are a bio-inspired notion of cloud-native applications. GLAs live in the cloud, and are able to migrate between data centers and cloud providers automatically, based on changes in cost and performance of cloud offerings, changes in customer behavior or requirements, or other factors.

The bio-inspired terminology applies for the different levels of components of a GLA:

- **GLA:** The GLA itself is a collection of organs that form the whole application.
- **Organ:** An organ consists of one or more cells.

- **Cell:** A cell is the lowest-level component of a GLA. There are two kinds of cells: The ones that provide user-defined functionality of the organ, and platform-defined cells like reverse-proxies for exposed ports

1.2 GLAPP

In order to manage these GLAs, we introduce a platform called GLAPP (**G**lobal **L**iving **A**pplication **P**latform). It allows a developer to deploy multiple GLAs on whatever cloud he/she has access to and sets a centralized mechanism in place to constantly monitor and manage all the GLAs. The platform supports heterogeneous environments, so a GLA can live and move across different providers, regions, instance types, etc.

GLAPP is an open-source platform which can be deployed easily on a server. Set-up instructions can be found on github¹, as well as all the source code that was developed in the course of this project²

1.3 Outline

This report serves the purpose of giving an overview over the whole project. In chapter 2, we give some background information about some design and implementation decision as well as the architecture of GLAPP. Chapter 3 introduces all the components GLAPP consists of in more detail. In chapter 4, we present the case study we had in mind and demoed in our final presentation. Chapter 5 discusses possible extensions and future work on the project. The report is concluded in chapter 6.

¹<https://github.com/glapp/gla-sails>

²<https://github.com/glapp>

6 PAGES: Background & Architecture

2.1 Basic Design Decisions

TODO ▷ *IN CHARGE: Adrian*◁

2.1.1 Main Components: Provisioning Backend, Frontend, MAPE

TODO ▷ *Describe the design decisions of the main components*◁

2.1.2 Deployment: Containerization with Docker

Containerization vs. other virtualization methods:

TODO ▷ *Explain the advantage of containerization compared to e.g. virtual machines*◁

Docker vs. other containerization implementations:

TODO ▷ *Explain the advantage of Docker compared to e.g. OpenVZ*◁

2.1.3 Orchestration: Docker Swarm

TODO ▷ *Explain the advantage of Docker Swarm compared to e.g. Kubernetes*◁

2.1.4 Rule-Based Adaptations vs. Markov Decision Process

TODO ▷ *DONE BY: Adrian / Riccardo*◁ **TODO** ▷ *Explain the thought process behind this decision*◁

2.2 Implementation Decisions

2.2.1 Backend: Node.js

The backend has been implemented with the *Node.js*¹ runtime. More specifically, we used a framework called *Sails*². Although it is labeled as a complete MVC framework, we mainly used

¹<https://nodejs.org/>

²<http://sailsjs.org/>

it to implement the server-side functionality and to provide a clean REST for the other GLAPP components.

There are two reasons why we chose Node.js and Sails in particular. For one thing, it was the backend framework some of us were most experienced in, which made it easier to focus on the actual project without having to work our way into a new framework. For another thing, Node.js comes with the well-established packaging system *npm*³, where we found very helpful packages for our problem domain. As an example, *dockerode*⁴ was the module which we used to implement the Docker-specific functionality, and it proved to be very well maintained, adding new functionalities whenever a new Docker release was published. Since we were about to work with relatively new technologies like Docker, it seemed reasonable to count on packages which we knew will be adapted to changes in these new technologies.

2.2.2 GUI: AngularJS

TODO ▷ *DONE BY: Dinesh* **TODO** ▷ *Explain the choice and alternatives*

2.2.3 MAPE: Java

TODO ▷ *DONE BY: Adrian / Riccardo* **TODO** ▷ *Explain the choice and alternatives* Java is a widely used programming language with good interoperability and decent support by other system or libraries used in this project. MAPE is a key component of the GLAPP platform. One of the important aspect of MAPE is the ability to easily interface with other components and system inside or outside GLAPP platform. Another consideration is the availability of third party libraries that GLAPP can leverage to perform common computation using well-known algorithm such as various planning and learning algorithms used in Markov Decision Process (MDP).

A set of comprehensive functionality is provided through Java API and third party libraries. For instance, though Java standard API, robust interfacing with SAILS backend monitoring system through HTTP connection is readily available. In addition, GSON library **TODO** ▷ *add reference* provides easy API to create and manipulate JSON objects that is used for data exchange between MAPE and SAILS backend as well as between MAPE and monitoring system. Most importantly, BURLAP library **TODO** ▷ *add reference* provides not only a set of planning and learning algorithm for Markov Decision Process in reinforcement learning, but also a framework for further extending the processing capability through custom implementation of various components including learning algorithm and approximation function. Availability of these functionalities from built-in and third party library makes Java a compelling language in developing MAPE component.

2.3 Architecture

The platform consists of 3 different parts/blocks: a frontend, a server backend and a control loop. The frontend provides an interface for developer to interact with the middleware to deploy and manage his/her GLAs. The server backend provides the management functionalities of the middleware including cloud infrastructure management, application deployment and access to the application status information. Lastly the control loop is the component responsible for enabling the management of GLAs by the platform itself. It follows the MAPE (Monitoring, Analysis, Planning and Execution) principle. Possible execution actions are moving cells between different cloud instances (migration), duplicating/splitting cells of the GLA (mitosis), or removing cells.

TODO ▷ *Explain the separate parts - DONE BY: Dinesh*

³<https://www.npmjs.com/>

⁴<https://www.npmjs.com/package/dockerode>

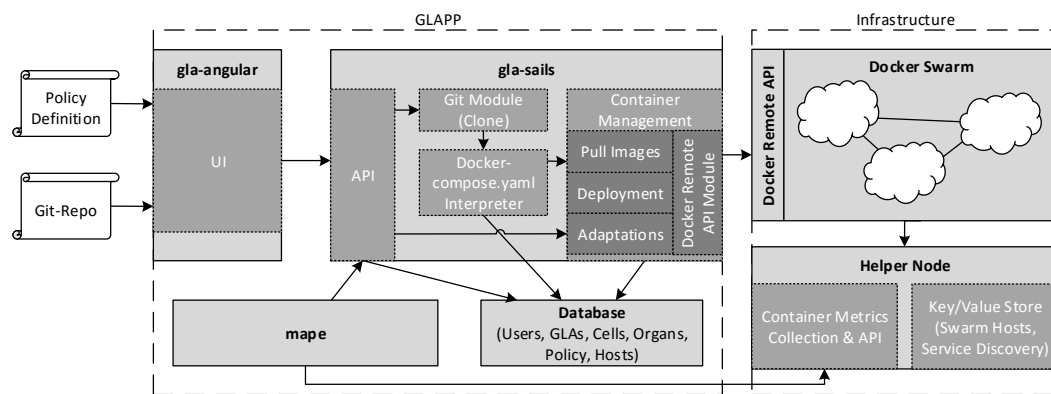


Figure 2.1: Detailed Architecture of GLAPP and the Infrastructure

10 PAGES: Components

3.1 Backend

The backend is mainly responsible for the deployment and is the central component that defines the data structure of the platform. In order to understand the main functionality of the backend server, this section introduces the data models, the REST API, the services and the tests of this component.

3.1.1 Data Models

Figure 3.1 provides a simple ER-like overview over the way how the server stores the data about users, applications, etc. The mentioned attributes are exemplary, but represent the most important ones to understand the whole model. It is apparent that the model is tightly coupled to the introduced notion of a GLA - an application which consists of organs, which again consist of cells.

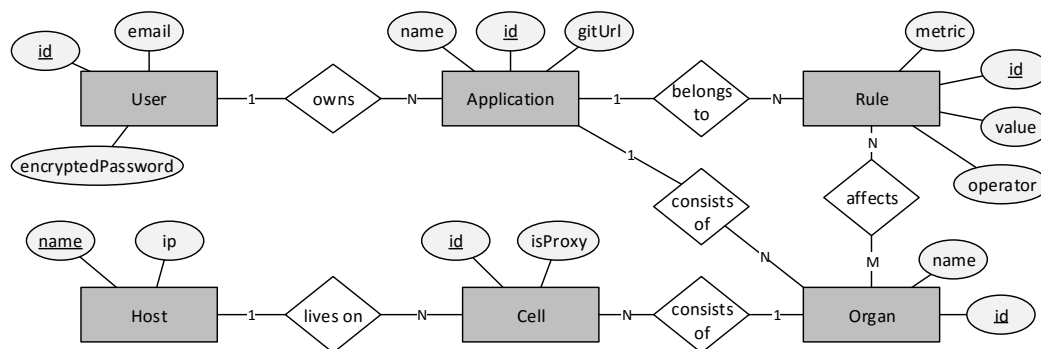


Figure 3.1: Simple overview over the models and their relations

Sails is built upon an ORM called *waterline*¹, which allowed us to manage the data on a relatively high level. That means that by defining the models and their relations and by setting up a connection to any database in the configuration of Sails, the data could be stored and retrieved by simple function calls. We decided to use *MongoDB*² as a database.

¹<https://github.com/balderdashy/waterline>

²<https://www.mongodb.com/community>

3.1.2 REST API

The backend is supposed to serve the frontend and the MAPE component with needed data and to provide an interface for triggering adaptations to the deployment. To achieve this clean interface, there is a REST API in place. Table 3.1 lists and describes the endpoints of this API. Since our components are completely modularized and the communication only happens via this interface, the remaining components like the frontend or the MAPE are easily interchangeable.

3.1.3 Services

In the Sails framework, services are used to handle specific, API-independent functionalities. For most of the internal logic apart from the actual API of the backend server, we created the following two well-encapsulated services:

- **DockerService:**

The DockerService is responsible for all the low-level, docker-specific functionalities. Many of the requests mentioned in table 3.1 trigger quite complex actions. For example, when a user requests to deploy an application, the server generally handles the following steps: (1) Creating an overlay network on the hosts, (2) Creating all the necessary cells in this network, and (3) Creating a reverse proxy cell for each organ that has an exposed port.

To avoid such low-level logic in the API controllers, it makes sense to encapsulate this in a service. The DockerService handles the connection to the Docker infrastructure by the mentioned npm package *dockerode*. It takes care of everything that is concerned with Docker, consisting of the following categories:

- Initialization: Establishing the Docker infrastructure connection, extracting the organs from the initial docker-compose.yml file
- Deployment: Handle the overlay network, deployment (creating the containers and the reverse proxies as cells)
- Adaptations: Moving containers, adding/removing containers, clearing all containers of a certain application
- Infrastructure: Get host information and parse it in order to save it in our database

- **PrometheusService:**

The PrometheusService is responsible for a clean provisioning of analytic data for the frontend. It establishes the connection to the Prometheus API and allows the fetching of any data the user requests. Since this data is cell-level data, it makes sense to have aggregation functions in place. In the current version of this service, only *average* is implemented, but it's programmed in a pipeline-like way so various alterations and aggregations of the data can be added and chained to the service call.

3.1.4 Tests

It is crucial for the whole platform that the backend's interface works as expected. To assure this, we wrote tests of all the API endpoints in different scenarios, based on the Node.js testing framework Mocha³. Since most requests trigger functions of the services, these tests cover the majority of code in the backend. How to run the tests is described in the README.md⁴ of the backend component.

³<https://mochajs.org/>

⁴<https://github.com/glapp/gla-sails>

Table 3.1: Endpoints of the REST API

Endpoint	Description
User	
<i>POST /user/signup</i>	Used by the frontend to sign up a user.
<i>PUT /user/login</i>	Used by the frontend to log in a user.
<i>GET /user/logout</i>	Used by the frontend to log out a user.
<i>GET /user/confirm-login</i>	Helper endpoint needed for the frontend to assure a user stays logged in on a refresh of the page.
Application	
<i>GET /application/getUserApps</i>	Used by the frontend to get high-level information about applications of a specific user.
<i>GET /application/details</i>	Used by the frontend to get detailed information of a specific application.
<i>POST /application/add</i>	Used by the frontend to add an application.
<i>POST /application/remove</i>	Used by the frontend to remove an application, including clearing up its existing cells on the hosts.
<i>POST /application/deploy</i>	Endpoint used by the frontend to deploy an application on the infrastructure.
<i>POST /application/undeploy</i>	Used by the frontend to clear up an application's existing cells on the hosts.
<i>POST /application/rename</i>	Used by the frontend to rename an application.
Organ	
<i>POST /organ/scaleUp</i>	Used by the MAPE to scale up an organ. Constraints can be defined to place the new cell on a specific host, e.g. in the region "US".
<i>POST /organ/scaleDown</i>	Used by the MAPE to scale down an organ. A specific cell can be specified to control what cell is going to be destroyed.
Cell	
<i>POST /cell/move</i>	Used by the MAPE to move a specific cell from one host to another.
Host	
<i>GET /host/info</i>	Used by the MAPE and the frontend get information about the current infrastructure.
<i>GET /host/prometheusUrl</i>	Used by the MAPE get obtain the URL of Prometheus.
Rule	
<i>GET /policy</i>	Used by the MAPE and the frontend to get the list of rules for a specific application.
<i>POST /policy/set</i>	Used by the frontend to set a new or replace an existing rule.
<i>POST /policy/remove</i>	Used by the frontend to remove a rule.
Analytics	
<i>GET /analytics/organCpu</i>	Used by the frontend to retrieve organ-level data for the cpu usage graph.
<i>GET /analytics/organMemory</i>	Used by the frontend to retrieve organ-level data for the memory usage graph.
<i>GET /analytics/events</i>	Used by the frontend to retrieve application-level data for the events graph.

3.2 2 PAGES: GUI

TODO ▷*DONE BY: Dinesh*◁

3.3 4 PAGES: MAPE

TODO ▷*DONE BY: Adrian / Riccardo*◁

MAPE is the control loop of the platform and responsible for managing the deployment of applications in accordance to the policy. It analyze the environment information and performance metrics to determine the healthiness of application and trigger appropriate adaptation action should there be any violation to the policy. The environment information includes details of the infrastructure, application deployment and user-defined policy.

MAPE consists of three parts: one is interfacing systems to retrieve the necessary information for analysis and dispatch the adaptation action, another part is analyzing the application healthiness in regard to the policy, and the last part is determining the adaptation action based on the environment information and computed healthiness.

At the start of each control loop, MAPE interfaces with SAILS backend to retrieve environment information, and with monitoring system to retrieve performance metrics. On the other hand, MAPE communicates to SAILS backend to execute the adaptation action. As mentioned in the implementation decision section, MAPE use HTTP protocol with JSON object for data exchange. It retrieves the infrastructure, application deployment information and user-defined policy from SAILS backend as well as performance metrics through HTTP requests. For instance, infrastructure information includes information of virtual machines such as the service provider, machine location and machine tier. Deployment information, on the other hand, refers to the information of each application components including on which containers they are currently deployed. Furthermore, MAPE also retrieve user-defined policy through SAILS backend. A policy is set of rules that defines which metrics the application need to comply to, its threshold for healthiness determine and the weight of each rule relative to the others. Lastly, MAPE retrieve the performance metrics from the backend monitoring system. JSON objects are used for all these information retrieval. The use of JSON objects provides a data structure that is flexible to represent different types of information from different interfacing systems. In addition, it allows MAPE to build a object-oriented view of the whole environment that provide quick access to specific information in various computation during the analysis.

TODO ▷*elaborate about healthiness computation*◁

After the healthiness values of the each individual cell is computed based on the policy and performance metric, all violating cells within an organ are taken into account for an evaluation. When the number of violating cell exceed a defined threshold, the organ will be considered as unhealthy. Once an organ is unhealthy, information of all violations will be further processed in MAPE to determine appropriate adaption action.

In the first implementation of MAPE, MDP is used to determine the adaptation action. Before starting the MDP, the environment and application deployment need to be modelled as states. At the same, actions for state transition are defined. BURLAP library and its framework is used for modelling and transition action definition. When the state and state-action information is constructed, learning algorithm from the library is applied. The learning algorithm is first applied in a simulated environment to verify correctness of the constructed model and defined transition action. It also serve a second purpose for understanding MDP process by generating the log of MDP with the aim for fine tuning the parameters in later stage of development. While MDP works close to the expected outcome in a simulated environment, concern was raised regarding the convergence of Q-value that learning algorithm rely on to determine the optimal action in a given state. More importantly, problem of high number of states and transition actions affects the viability of applying MDP to the optimization problem of MAPE. After subsequent iteration of adjustment and refinement to the model and transition actions and execution of MDP in the simulated environment, it is concluded that reduction of

number of states is required for applying MDP. To reduce the number of states, value function approximation need to be used.

3.4 2 PAGES: Other

3.4.1 Monitoring / Prometheus

TODO ▷*DONE BY: Adrian / Riccardo*◁

Prometheus **TODO** ▷*add reference*◁ is a monitoring system that can provide both infrastructure and application performance metrics. Infrastructure metrics are provided by exporters, which is a set of library that expose metrics data of the host environments and the Docker containers. Application metrics are provided with the use of custom-built program to expose the needed data in a Prometheus exposition format **TODO** ▷*add reference*◁. The use of Prometheus and its exposition format enable MAPE to access custom metrics such as application performance data or cost metrics of the host from various cloud service providers.

Prometheus is also sophisticated in providing different types of query. It supports range query and point query with additional functions that can be applied on the raw data when formulating a query. This enables the fine-grained extraction of data that best fits the corresponding performance metric. For instance, the rate function allows the direct retrieval of the rate of change in CPU time of a container with configurable data point range, data point interval and duration covered by each data point. The robust query functionality provides two major advantages. First of all, it reduces the complexity of computation logic of MAPE that is responsible for getting specific performance metric for subsequent healthiness computation. Secondly, it allows quick tweaking of query parameters to obtain the most relevant data for adaptation action determination during later development stage.

3.4.2 Supporting Tools & Applications

TODO ▷*Infrastructure Scripts, Consul, Registrator, Voting-App, Proxy, metrics-server etc.*◁ In the course of the project, several further small tools and applications were used to complete the task. This section quickly mentions them with a reference to further information sources:

- **Infrastructure scripts:**

To always work on the same host infrastructure, we soon created handy scripts to set up a complete infrastructure which includes all of the required components. This is all collected in a git repository⁵ and can also be used by other people to set up the same infrastructure as we were working with.

- **Consul.io:**

As a discovery service for the Docker swarm hosts, we used the key/value store of Consul.io⁶. It is included in the mentioned infrastructure scripts.

- **Registrator:**

To register every exposed port of the deployed applications, we used Consul.io in combination with Registrator⁷. In order to register not only public ports but also exposed ports that are only available in the created overlay network, the *-internal* run option proved to be very helpful.

- **HAProxy:**

Registrator was used to create the mentioned reverse proxies. In order to always have an up-to-date reverse proxy with a complete list of active endpoints of a specific organ, we

⁵<https://github.com/glapp/docker-swarm-creation>

⁶<https://www.consul.io/>

⁷<http://gliderlabs.com/registrator/latest/>

created a HAProxy docker image⁸ that used Registrator and consul-template⁹ to dynamically adjust the configuration of the HAProxy and restart it with every change.

- **Voting-App:**

The voting-app is our demo application that will be introduced in more detail in the chapter 4. It is based on an example application by Docker itself¹⁰. We adjusted the individual components to our needs, e.g. to include custom metrics, and they can all be found on our github page.

- **Metrics-Server:**

The metrics-server¹¹ is a small service we created that was used to register the clicks of the demo application. It could be extended to handle all different kinds of metrics via a REST API.

⁸<https://github.com/Clabfabs/miniboxes/tree/master/minihaproxy>, <https://hub.docker.com/r/clabs/haproxy/b/>

⁹<https://github.com/hashicorp/consul-template>

¹⁰<https://blog.docker.com/2015/11/docker-toolbox-compose/>

¹¹<https://github.com/glapp/metrics-server>

7 PAGES: Case Study

TODO ▷ *DONE BY: Riccardo* ◁

4.1 Explain Demo Application

- Explain concept of the app at high level
- Explain the components themselves
- Explain the docker-compose file
- Explain our custom metrics (costs, clicks)
- Explain modifications to the app (region switch button, POST request with click-origin to metrics-server with every click)

4.2 Scenarios

- Explain the scenarios
- Explain the triggers
- Explain how MAPE realizes a necessary adaptation

2 PAGES: Future Work

TODO ▷ *DONE BY: Fabio / Adrian*◁

1 PAGE: Conclusion

Bibliography