

Design and implementation of a social network infrastructure for designers of Multi-Cloud applications

Christos Papoulas

Thesis submitted in partial fulfillment of the requirements for the

Masters' of Science degree in Computer Science

University of Crete

School of Sciences and Engineering

Computer Science Department

Knossou Av., P.O. Box 2208, Heraklion, GR-71409, Greece

Thesis Advisors: Prof. *Kostas*, Dr. *Magoutis*

UNIVERSITY OF CRETE
COMPUTER SCIENCE DEPARTMENT

Your Title

Thesis submitted by
Author Name
in partial fulfillment of the requirements for the
Masters' of Science degree in Computer Science

THESIS APPROVAL

Author: _____
Author Name

Committee approvals: _____
Name of first member
Assistant Professor, Thesis Supervisor

Name of second member
Associate Professor, Committee Member

Name of third member
Professor, Committee Member

Departmental approval: _____
Name of Director of Graduate Studies
Professor, Director of Graduate Studies

Heraklion, July 2015

Abstract

In this work we propose a Social Network Platform, both front-end and back-end technology, for model-driven software design and deployment of applications in Multi-cloud environments. Nowadays, DevOps users, especially cloud deployment specialists, wander around the web looking for automated tools like Chef supermarket, IBM Bluemix and other deployment tools in order to, almost manually, configure and deploy their applications without the assist of any community-sourced information repository.

The innovation of our Social Network Platform is that its users can create their own Cloud Application Modelling and Execution Language models or benefit from automatically generated models. We provide a Platform with an integrated community assisting the DevOps users design and deploy their applications and a repository of execution histories of applications. We incentivize user to stay inside the Platform instead of roaming around the web looking for other Q&A sites, such as StackOverflow, to find their answers. The information about the runtime executions of applications can be uploaded and analysed by the Platform presenting cost effectiveness analysis. In addition, Natural language Processing tools are integrated with the Platform in order to identify the context of users questions and provide automated answers. A scalable front-end is implemented using multiple front-end nodes and a powerful back-end system is implemented using caching technologies.

Περίληψη

Στην εργασία αυτή παρουσιάζεται μία πλατφόρμα κοινωνικής δικτύωσης, η τεχνολογία front-end καθώς και η τεχνολογία back-end αυτής της Πλατφόρμας. Η πλατφόρμα αυτή προορίζεται για την σχεδίαση εφαρμογών οδηγούμενων από μοντέλα και την εγκατάσταση των εφαρμογών αυτών σε περιβάλλοντα πολλαπλών νεφών. Σήμερα, οι χρήστες για την Ανάπτυξη και Λειτουργία Εφαρμογών, ειδικά αυτοί που ειδικεύονται στην ανάπτυξη εφαρμογών σε νέφη, περιπλανιούνται στο παγκόσμιο ιστό αναζητώντας αυτόματα εργαλεία όπως το Chef supermarket, το IBM Bluemix και άλλα όπου σχεδόν χειροκίνητα προσαρμόζουν και εγκαθιστούν τις εφαρμογές τους χωρίς την βοήθεια από κάποια κοινότητα που να παρέχει πληροφορίες και να αποθηκεύονται σε κάποιο χώρο δεδομένων. Οι χρήστες, μέσα από την προτεινόμενη πλατφόρμα, μπορούν να δημιουργήσουν τα δικά τους Cloud Application Modelling and Execution Language μοντέλα ή να επωφεληθούν από αυτόματα παραγόμενα μοντέλα. Παρέχουμε μια πλατφόρμα με μια ολοκληρωμένη κοινότητα που βοηθάει τους χρήστες να σχεδιάσουν και να εγκαταστήσουν τις εφαρμογές τους, καθώς επίσης και ένα χώρο δεδομένων από προηγούμενες εκτελέσεις των εφαρμογών. Παροτρύνουμε τους χρήστες να παραμείνουν μέσα στην Πλατφόρμα από το να περιπλανούνται στο παγκόσμιο ιστό ψάχνοντας σε άλλες ιστοσελίδες με Απαντήσεις & ερωτήσεις όπως το StackOverflow ψάχνοντας να βρουν τις απαντήσεις τους. Η πληροφορία για τα περιβάλλοντα εκτέλεσης των εφαρμογών μπορεί να μεταφορτωθεί και να αναλυθεί από την Πλατφόρμα παρουσιάζοντας την ανάλυση της αποτελεσματικότητας του κόστους. Επιπροσθέτως, εργαλεία για Επεξεργασία Φυσικής Γλώσσας έχουν προστεθεί στην Πλατφόρμα με σκοπό την αναγνώριση του εννοιολογικού περιεχομένου των χρηστών που κάνουν ερωτήσεις ώστε να παρέχονται αυτοματοποιημένες απαντήσεις. Ένα Κλιμακώσιμο front-end υλοποιήθηκε χρησιμοποιώντας πολλαπλούς front-end κόμβους και επίσης υλοποιήθηκε ένα ισχυρό back-end σύστημα χρησιμοποιώντας τεχνικές προσωρινής αποθήκευσης.

Acknowledgements

Test acks ...

Contents

1	Introduction	1
1.1	Background	2
2	Related Work	5
2.1	Caching Application Data	5
2.2	Professional Networks	6
2.3	Natural Language Processing	9
2.4	Configuration management and deployment	10
3	Implementation	11
3.1	Implementation of Social Network	12
3.1.1	CDO communication with Social Networking Platform	16
3.2	Scaling Social Network Engine	17
3.3	Memcache	18
3.4	Natural Language Processing and classification	19
3.4.1	Bayes Classification Algorithm	22
3.4.2	Automated answers from Natural Language Processing	22
3.5	User interface	23
3.5.1	User Interface Design Principles	23
3.5.2	Automated Application Model Creation	25
3.5.3	Graphical modeling of applications	27
4	Evaluation	29
4.1	Improving Performance with memcached	30
4.2	Improving Performance with engine	32
4.3	Evaluation of NLP classification	32
4.4	Evaluation of requirements and UI interface	33
5	Conclusions and Future Work	37

List of Figures

1.1	CAMEL DSLs.	3
3.1	The overall architecture of Social Network.	11
3.2	Architecture of the Elgg Social Networking engine.	12
3.3	The Elgg Engine Data model.	13
3.4	The structure of the application description plug-in.	15
3.5	The scenario a depicts a request from memcached when the key does not exist and scenario b depicts a updated operation of a value)	19
3.6	The main StackOverflow users' actions and NLP Classifier.	21
3.7	Automated answer to user's question using NLP.	23
3.8	The engineering & social activities are seamlessly within the Platform.	23
3.9	The application model home page	24
3.10	Steps for Automated creation of baseline model.	26
3.11	Final Step of Automated creation of baseline model.	26
3.12	GMF editor composition of a sample application.	27
4.1	The average response time for all configurations.	31
4.2	The average CPU utilization for all components.	34
4.3	The Response time for two Social Network Engines.	35
4.4	The CPU utilization for two Social Network Engines.	35

List of Tables

2.1	Feature comparison.	8
4.1	Number of Queries to Social Network and CDO server Databases. .	31
4.2	VM resources	32
4.3	NLP Evaluation of Classification	33

Chapter 1

Introduction

In this work, the design and implementation of a social networking platform for designers of Multi-Cloud applications is presented. In this targeted social networking platform, DevOps engineers [1] and particularly cloud deployment specialists can benefit from other users' experience and answer design questions such as which is the most cost-effective deployment and which configuration best fits their needs.

DevOps community consists of several types of specialists from both Development and Operations fields, such as programmers, testers, Quality Assurance staff [2] from the first and system administrators from the second. All the above professionals require different tools for various purposes like building applications, testing, deploying routines, configuration, automation of utilities, tracking and versioning in systems. Cloud Deployment Specialists are responsible for migration and configuration of hosted solutions for applications. Furthermore, they have deep industry knowledge about security, auto scaling, storage, load balancers, Content Delivery Networks [3] and everything related to an application's secured hosting and scalability. The communication between these users is essential in order to exchange opinions and solutions on cloud deployment of applications. Having a repository, where they can find the configurations and the executions data of an application and point to those data is valuable since it can make their conversation specific and help them come to more concrete conclusions.

This social network targets the cloud deployment specialists and binds all social networking concepts such as personal messaging, groups, new feeds with concepts of application composition and deployment, integrating a repository of cloud applications and infrastructure description based on Cloud Application Modelling and Execution Language (CAMEL) [4].

This CAMEL repository brings to several benefits to the DevOps community. Among the range of possible DevOps tasks, the Social Network focuses on selecting the most appropriate deployment configuration for an application. This is especially challenging in a multi-cloud setting due to the large diversity of deployment possibilities and tradeoffs.

Currently DevOps users work with a small set of well-understood deployment

options, missing on opportunities for improving performance, reliability and/or lower cost. Investigation of new options involves time consuming testing over new infrastructures. Discussing those topics with the community in online social or technical forums may provide insight over deployment options; however, the answer to a hard question often needs to be backed by experimental data that is not readily available.

An integrated environment like our proposal, comes to solve the above issues by enriching user interactions with structured references to applications and their components, execution data, and mined knowledge from real deployments. Mined knowledge can be combined with user activity and profiles to provide personalized suggestions and hints. An improved mode of user interaction is expected to result in stronger incentives for DevOps users to contribute information to the underlying repositories. Richer content should lead to better quality of mined knowledge, benefiting the DevOps community and providing further incentive for contributions. The social networking platform designed to be closely integrated with a set of information repositories satisfying the following requirements: (R1) handle entire applications rather than just software components; (R2) abstract application structure through software modeling; (R3) capture and analyze application runtime performance.

Altering the focus from the individual components of the applications to the whole applications and the analysis of its execution data, can provide answers to many interesting questions of the community and support discussions and arguments with hard data. These requirements can provide software developers with strong urge to contribute, leading to the sustainability and growth of information and derived knowledge in the repository.

This thesis is structured as follows: in the section 1.1 the background of application models and the CAMEL repository in the context of PaaSage EU project is presented. Chapter 2 describes the previous work based on Caching architectures, other Social Networking Platforms and the Natural Language Processing. In the chapter 3, the system implementation of the Social Network is described and the chapter 4 evaluates the implemented system. Finally, the present thesis concluded in chapter 5.

1.1 Background

This section presents an overview of the descriptions of the application models that are available on the Social Network Platform and a summary of the technologies that are used by the PaaSage in order to assist the cloud deployment specialists to deploy the aforementioned application models.

Application models inside social network platform are described in CAMEL. CAMEL integrates various domain-specific languages (DSLs). DSLs provide a notation tailored towards an application domain and are based on the relevant concepts and features of that domain. As such, a DSL is a mean to describe

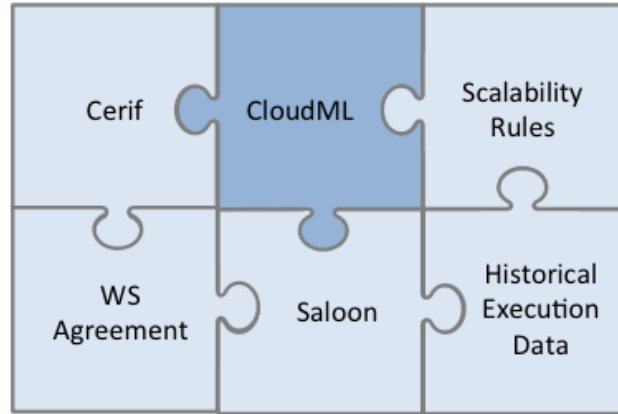
and generate members of a family of programs in the domain. These DSLs cover a wealth of aspects of specification and execution of multi-cloud applications like CloudML, Scalability rules, WS Agreement, Saloon and Historical Execution Data.

CloudML [5] is a recent approach that focuses on the provisioning and deployment of multi-cloud applications, is built upon MDE techniques and methods, and provides a models@run-time [6] environment for enacting the provisioning. WS Agreement [7] is a Web Services protocol for establishing agreement between two parties, such as between a service provider and customer. Saloon [8] is an approach that uses models to represent clouds variability, as well as ontologies to describe the heterogeneous aspect of the cloud ecosystem. The CAMEL model assembles all those DSLs as shown in figure 1.1.

CAMEL is using the Eclipse Modelling Framework (EMF) [9] on top of the Connected Data Objects (CDO) [10] [11]. Application Models are persisted on the CDO repository. EMF is used as a building tool for the application models. CAMEL model specification is written in XMI format (XML Metadata Interchange) which is a standard for exchanging metadata information via Extensible Markup Language (XML). Thus, XMI consists the main describing system of CAMEL and EMF provides tools that enable viewing and command-based or tree-based editing of the application models.

The CAMEL repository is currently being populated with a wealth of information from multi-cloud deployments of various distributed applications [12]. Also, SNP performs analytics over the CAMEL repository to extract knowledge about deployments characteristics that work best for certain applications and use it in the context of the professional network.

Figure 1.1: CAMEL DSLs.



The Social Network Platform that this thesis proposes is part of the PaaSage EU Project [13]. The PaaSage perspective is to be a tool for a cloud deployment Specialist to leverage the complex task of deploying an application to the clouds. Usually, a cloud deployment specialist could easily learn the interface and features of one Cloud provider, but it would be very costly and time consuming to lever-

age the development to many providers. It is a real challenge to orchestrate the simultaneous deployment to many different Clouds at the same time. The main objective of PaaSage is to assist the developer to deal with difficult deployment scenarios through automatic cloud deployment. In order to satisfy this, several components are included to the PaaSage ecosystem. The Profiler components read the CAMEL models and convert them into a constraint programming model by defining the variables of the model, their domains, and the constraints that must be satisfied by the deployment. Also, the Profiler checks all constraints of the CAMEL model and sets the domains of the variables accordingly. The reasoning component analyses the model and finds how deployment candidates should be evaluated. Once a solution has been found, the reasoning component converts the model to Cloud Provider Specific Models (CPSM) for the providers involved in the proposed deployment. The adapter component takes the CPSMs, produces and validates a configuration plan, and sends this plan to the execution ware. The execution ware [14] receives the deployment plan from the adapter and enacts the deployment of the application on the selected providers. Furthermore, the execution ware interacts with the Cloud providers, acquires the virtual machines, configures them and launches the user application on the set of virtual machines. Once the machines are running, the execution ware collects sensor data for the running application, triggering re-configurations if necessary.

The Social Network Platform brings to the DevOps users a friendly interface to browse, discover, view and discuss Application Models. Furthermore, it presents a way to deploy and run these Application Models, by using the previously mentioned components under the hood, and mines their execution history data.

The key contributions of this Social Network Platform(SNP) are the following:

- The SNP binds all the Social Networking aspects such as friends, new feeds, personal messages etc. with the engineering aspects of creating and deploying application models.
- The SNP brings the execution histories of the CAMEL applications in the light, providing the end users the ability to browse, discuss, point and find essential information needed for other applications.
- The SNP uses the best known practices both for the front end viewing system and the back end technology.
- The SNP runs on a horizontal scale architecture with memcached at the back end to reach near real time interaction.

Chapter 2

Related Work

In this chapter the related work of other professional networks and their caching architecture is described in section 2.1. In addition, an overview and related work of Natural Language Processing is presented in section 2.2.

2.1 Caching Application Data

Arguably two of the largest existing networking platforms are Linkedin and Facebook. The caching technologies of these networks are of great interest for the way they are managing and storing vast amounts of data.

Linkedin, the largest professional network, stores hundreds of terabytes of data to Project Voldemort [15], a key-value store, inspired by Amazon Dynamo [16], another well-known key-value store. Linkedin stores to Voldemort pre-computed offline data. For example, it stores the results of data mining applications, such as features like “People You May Know”, that are running on hundreds of terabytes to make an estimation and are using Hadoop as the computational component of those estimations. Voldemort and Dynamo have the same following requirements: (1) a simple *get/put* application interface (2) A *replication* factor, the number of replicas for each key-value tuple, implemented using vector clock, (3) a *required read* factor to succeed a get request and (4) a *required write* factor to succeed a put request.

Facebook, the largest social network, serves billions of requests per second using memcached [17]. In this magnitude of scale, Facebook has several pools of memcached servers (regional pools) around the globe. A request for a single page can produce hundred of requests to the back-end system. Memcached is used to store not only key-value from MySQL queries but also pre-computed results from sophisticated algorithms. In order to achieve a near real time communication experience to the end user, memcached servers have to be efficient, reducing latency to minimum.

The research question in such systems is when a particular key will be invalidated. This problem occurs according to [17] in two cases: (1) *stale sets* and (2)

thundering herds. A stale set occurs when a web server sets a value to the memcached that does not reflect the real value of the database. Thundering herds occur when a specific key has a heavy read and write activity at the same time. Stale sets are resolved by an N-bit token, that is bound to a specific key and sent from the memcached to the web server that wants to update the key when a cache miss occurs. If a delete request is received, the request for updating this value from the client is rejected. The thundering herds are solved by configuring the memcached servers to return an N-bit token only once every ten seconds per key.

The PaaSage Social Network, inspired by Facebook, integrates memcached in its back-end architecture and it is configured to properly interact with the caching application.

2.2 Professional Networks

This section reviews the related work on other professional social networking platforms and what they provide. Having the requirements of DevOps users as a priority, we compare those platforms with our proposed SNP, pointing their deficiencies and filling in the gaps with our innovation.

The table 2.1 summarizes and categorizes the characteristics of the most important related approaches along the following dimensions (depicted as columns of table 2.1):

- which of the following key social features are supported by the platform: follow users; news feeds; groups; Q&A; personal messages;
- does the platform rely on one or more repositories to store the following type of information: software code, software models, configuration information, execution histories; and whether these repositories are community sourced;
- does the social networking platform leverage the repositories to provide users with specific suggestions and hints;
- does it support application deployment?

Information Technology (IT) [18] professionals use a variety of online sources as aids in their daily tasks. Developers typically prefer community-moderated forums over vendor-moderated sites [19]. Social networks focusing on software technology in particular provide developers with the opportunity to leverage the knowledge and expertise of their peers.

One of the most popular such platforms is GitHub [20], a collaborative revision control platform for developers launched in April 2008, and arguably the largest code-hosting site in the world. GitHub provides social networking functionality such as feeds, followers, wikis and a social network graph that captures how developers work on versions of their repositories, which version is newest, etc. Gitter [21] is a related service that facilitates discussions between members

of GitHub communities by providing a long-term chat integrated with code and issues. Sourceforge [22] was the first code-hosting platform offered to open-source projects. It was launched in 1999 and offered IT professionals the ability to develop, download, review, and publish open-source software. Sourceforge is similar to GitHub in its support for social features. Other similar code-hosting platforms are Google Code [23] and Microsoft CodePlex [24]. None of those platforms collect, analyze, or use information from executions of application deployments to improve the level of technical discussion between users or abstract code structure through modelling or enhance user interactions through the use of analytics over application execution histories.

StackOverflow [25] advances on earlier community-driven Q&A sites in which users ask and answer questions. Users can vote up or down questions and answers and earn *reputation points* and *badges* in return for their active participation. Although StackOverflow and GitHub address different aspects of software development (StackOverflow is not a code-hosting platform) there is a synergy and correlation between the two [26]. The proposed social network platform extends StackOverflow through the use of social networking features that enable users interested in reasoning about application deployments to use and share knowledge drawn from analyses of information repositories.

IBM's BlueMix [27] is a development and support platform for communities of DevOps users wishing to compose distributed applications out of components drawn from libraries and deploy them at IBM-provided and supported cloud infrastructure. BlueMix is a key component of IBM's DevOps best practices [28] for achieving rapid prototyping, automated deployment, and continuous testing of software. BlueMix encourages its users to ask their questions to StackOverflow but also includes a community forum [27] with rating of answers contributes to eventually building a basic knowledge base, similar to traditional approaches such as StackOverflow. The proposed social network platform system differs from BlueMix in its support for expressing applications as models (CloudML, CAMEL) and its use of two information repositories, the PaaSage repository of models and execution histories and Chef supermarket, and the use of analytics over past executions to enable users to reason about application deployments. A common feature between the proposed social network platform and BlueMix is support for deployment of distributed applications.

Linkedin is widely adopted across a range of professional communities due to its robust set of social features (and to some extent due to its use of extensive analytics over collected information [29]), LinkedIn provides no specific support for software engineering activities and thus more closely resembles traditional social networking platforms such as Facebook.

The lack of Social Networking features of github came to fill the Geeklist platform [30], where developers and IT companies can discover and share the work they have done, connect with other companies in a social network manner or join development communities. Another code hosting platform is Snipplr [31], where developers can upload short code snippets but not full programs, in order to keep

Table 2.1: Feature comparison.

	User Interaction				Repository					Repo assisted hints ^b	Application deployment
	Social features ^a	Groups	Q & A	Personal messaging	Software code	Software models	Software config	Execution histories	Crowd sourced		
GitHub	✓	✗	✗	✗	✓	✗	✗	✗	✓	✗	✗
Sourceforge	✓	✗	✓	✓	✓	✗	✗	✗	✓	✗	✗
GoogleCode	✗	✗	✓	✗	✓	✗	✗	✗	✓	✗	✗
CodePlex	✓	✗	✓	✓	✓	✗	✗	✗	✓	✗	✗
StackOverflow	✗	✗	✓	✗	✗	✗	✗	✗	✗	✗	✗
BlueMix	✗	✗	✓	✗	✗	✗	✓	✗	✗	✗	✓
Chef Supermarket	✗	✗	✗	✗	✗	✗	✓	✗	✓	✗	✗
LinkedIn	✓	✓	✓	✓	✗	✗	✗	✗	✗	✗	✗
geeklist	✓	✓	✗	✗	✗	✗	✗	✗	✓	✗	✗
Snipplr	✗	✗	✗	✗	✓	✗	✗	✗	✗	✗	✗
Masterbranch	✗	✗	✗	✗	✓	✗	✗	✗	✗	✗	✗
Dzone	✓	✗	✗	✗	✓	✗	✗	✗	✗	✗	✗
codeproject	✓	✗	✓	✗	✓	✗	✗	✗	✓	✗	✗
PaaSage SN	✓	✓	✓	✓	✗	✓	✓	✓	✓	✓	✓

^a Features: follow and news feed^b User assistance based on data analysis of the repository

all of their frequently used code in one place that is accessible from any computer and any user. Masterbranch [32] is a new under development platform that allows collating and sharing of projects within a user's profile. This profile works similarly to LinkedIn and has an incentivisation scheme called DevScore, coupled with unlockable achievements that add a gamification element. Dzone [33] is essentially a link repository for developers allowing link sharing and incentivisation based on voting for the popular links. The Code project [34] website and forum allow code-specific discussion and share relevant articles and news, contains blogs, newsletter and a questions and answers section.

The above systems can be further classified based on whether they use a repository to store software-related information (code, models, configuration, or execution histories) and whether this information is shared and raised through crowd sourcing [35]. GitHub, GoogleCode, CodePlex, SourceForge, BlueMix, Chef supermarket, and our platform store at least one type of software-related information and all systems but BlueMix are raising shared content in their software-related repositories via crowd sourcing. Our Social Network Platform is the only solution that analyzes information in its software-related repositories in order to provide users with assisting suggestions and hints. The Platform targets the model of the application and the Software code of the application is not stored inside the Social Network Platform.

2.3 Natural Language Processing

Another technology that has been frequently used by networking platforms and we have integrated to our SNP is Natural Language Processing (NLP) [36]. Specifically, we focus on social networking usage of NLP and how knowledge can be mined from repositories of Q&A sites.

NLP has been used to process Twitter's messages and come to some results according to the classifications. Twitter has a good pool of micro-blog text which is suitable for NLP because of the small text sentences that users are allowed to post. Those posts describe emotions, feelings, opinions or situations. So several techniques [37] [38] [39] have been introduced to process and classify twitter posts in several categories.

StackOverflow (SO) is not left without NLP, because it can be seen as a repository of Q&A for programming questions. This means, that most of the questions and answers in SO contain some kind of a description at first and some code afterwards. An autoComment tool [40] is proposed using NLP which maps the code from developer projects and locates the same code somewhere in the SO Q&A, if it exists. If autoComment matches a segment of the developer's code with a code segment at SO, it performs NLP to the description of the code and inserts the modified description to the developer's code.

A trend in Social Networking sites is the ability of users to "tag" their posts. Those tags describe the users' goals and interests. Tagging SO questions involves

askers selecting appropriate keywords to broadly identify the domains to which their questions are related. There also exist mechanisms by which other users can subscribe to tags, search via tags, mark tags as favorites, etc. This users' classification of context is used by PaaSage Social Network Platform as described in section 3.4.

Social Network Platform uses those tags and Natural Language Processing to answer the following research questions:

- Can the platform identify the similarity of a given question with other questions already posted in the system.
- Can the platform map a question with a relevant query to the repository in order to provide the one who asks with an appropriate response.
- Can the platform paraphrase the queries according to the user's arbitrary input in order to meet the previous objective.

2.4 Configuration management and deployment

Another key component in a portfolio of DevOps tools is configuration management (CM) [41], the process of maintaining a detailed recording of software and hardware components in an infrastructure. An effective CM process provides significant benefits including reduced complexity through abstraction, greater flexibility, faster machine deployment, faster disaster recovery, etc. There are numerous configuration management tools from which a system administrator can choose, however the most widely known are: Bcfg2 [42], CFEngine [43], Chef [44] and Puppet [45]. Each of these tools has its strengths and weaknesses [46], [47]. In a DevOps environment, a CM solution is often combined with provisioning and deployment tooling [28].

The Social Network Platform uses Chef as a CM and deployment automation tool to support professional network users, and SNP integrates the Chef cookbooks in its platform.

Furthermore, a recent trend in DevOps software development is continuous integration (CI) [48] and automated code deployment and testing off of online code repositories. Travis [49] is a CI tool that automatically detects when a commit has been made and pushed to a GitHub repository, subsequently tries to build the project, deploy and run tests, and notify the user of the status. Another popular CI tool is Jenkins [50], an open-source software tool for testing and reporting on isolated code changes in real time. Similar to Travis, Jenkins enables developers to find and solve defects in their code rapidly and automates the testing of their builds.

Although the PaaSage social networking platform does not provide a complete CI solution, it automates the deployment of complex applications through a model-driven process.

Chapter 3

Implementation

This section describes the implementation of social network site, the User Interface and how the system scales.

The system is composed by the following components, as shown in figure 3.1: At the first layer lives (1) the Social Networking engine, which runs all PHP scripts and described in section 3.1. At the second layer lives (2) the Memcached caching system, which described in section 3.3. At the third layer lives (3) the Social Network MySQL database, and (4) the CDO server - client components and the CDO repository. The Social Networking Engine at layer 1 is considered as the front end system and the CDO Client, the memcached nodes, the CDO Server and the repositories are considered as the back end system.

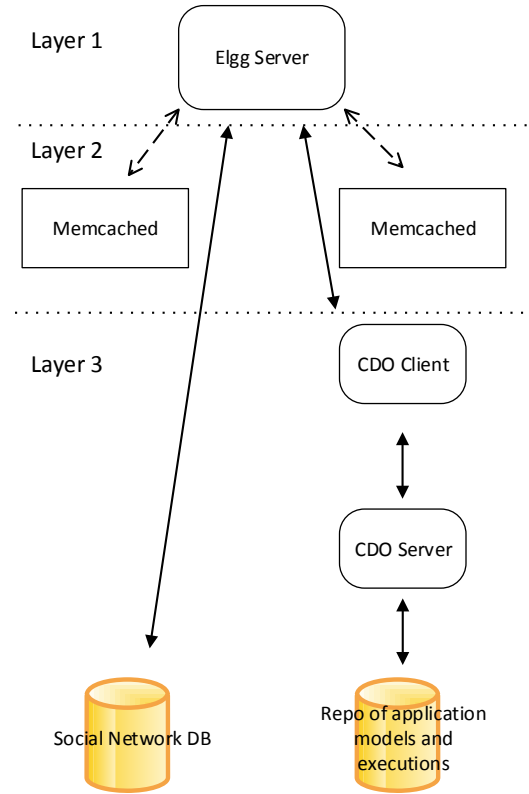
Achieving the scalability of the system, two system architectures are examined at two layers of the system: (1) We added more than one Social Network engine at the first layer of the system. In this implementation, in order to keep the file system in consistent mode we integrated Apache Zookeeper [51] as described in section 3.2. (2) We added more than one memcached nodes at the second layer in order to add more cpu capacity and improve the system response time as described in section 3.3.

3.1 Implementation of Social Network

The social networking platform is implemented over the extensible Elgg social network framework [52]. Elgg is open source software written in PHP, uses MySQL for data persistence and supports jQuery [53] for client-side scripting.

The jQuery is preferred instead of pure JavaScript [54] since it's a light library which pushes content to the client machine, it therefore reduces the wait time for server response. Plus, it's smaller than Flash, so it results in smoother playbacks and less errors. Furthermore, jQuery works anywhere since is cross-browser compatible with any browser, mobile phone or tablet, and Apple devices. Finally, jQuery's syntax is designed to make it easier to navigate to a document, select an HTML DOM elements, create animations, handle events, and developing Ajax

Figure 3.1: The overall architecture of Social Network.

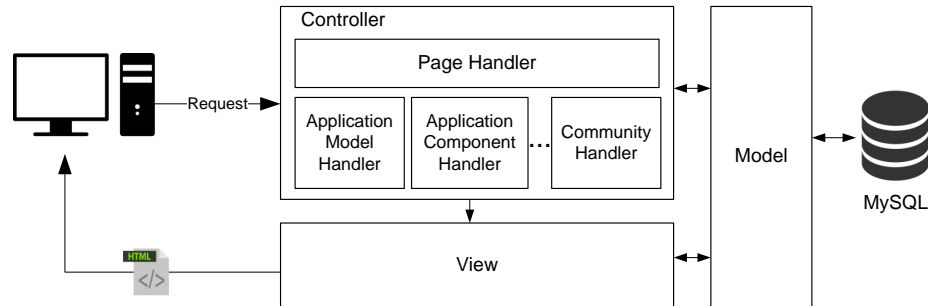


applications. So, the jQuery is used for implementing client side scripting and in the remaining of this chapter, when we refer to JavaScript, we actually refer to the jQuery library. Furthermore, some other JavaScript libraries are used in order to make the User Interface more powerful. The Chart.js [55] library is used to generate the graphs and charts in execution's page.

The architecture of Elgg Social Network shown in figure 3.2. The Model of the framework is structured around the following key concepts as shown in figure 3.3:

- *Entities*, classes capturing social networking concepts: users, communities, application models. Elgg Core comes with four basic objects: ElggObject, ElggUser, ElggGroup, ElggSite, ElggSession, ElggCache and a lot of other classes necessary for the proper engine operation.
- *Metadata* describing and extending entities (e.g., a response to a question, a review of an application model, etc.).
- *Relationships* connect two entities (e.g., user A is a friend of user B, user C

Figure 3.2: Architecture of the Elgg Social Networking engine.



is a contributor to an application model, etc.) and are persisted in the Social Network DB.

- *Annotations* are pieces of simple data attached to an entity that allow users to leave ratings, or other relevant feedback.

All Elgg objects inherit from *ElggEntity*, which provides the general attributes of an object. Elgg core comes with the following basic entities: *ElggObject*, *ElggUser*, *ElggGroup*, *ElggSite*, *ElggSession*, *ElggCache*, as well as other classes necessary for the operation of the engine.

Figure 3.3: The Elgg Engine Data model.

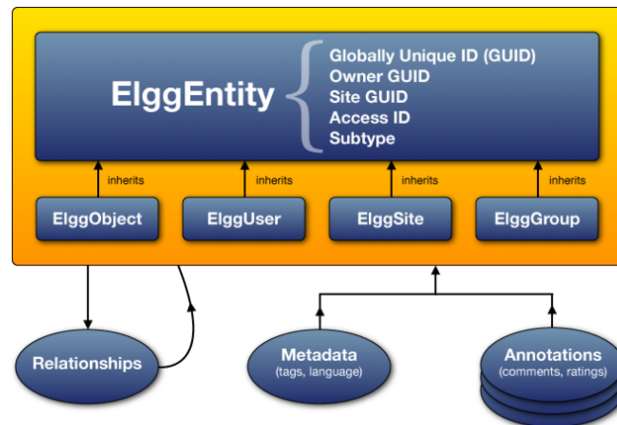


Figure 3.2 shows the model, view, and control parts of Elgg's architecture. In a typical scenario, a web client requests an HTML page (e.g., the description of an application model). The request arrives at the *Controller*, which confirms that the application exists and instructs *Model* to increase the view counter on the application model object. The controller dispatches the request to the appropriate handler (e.g., application model, component handler, community handler) which

then turns the request to the view system. View pulls the information about the application model and creates the HTML page returned to the web client.

Elgg comprises a core system that can be extended through plugins (examples are the Cart system or the handling of Application Models). Plugins add new functionality, can customize aspects of the Elgg engine, or change the representation of pages. A plugin can create new objects (e.g., `ApplicationObject`) characterized (through inheritance of `ElggEntity`) by a numeric globally unique identifier (GUID), owner GUID, Access ID. Access ID encodes permissions ensuring that when a page requests data it does not touch data the current user does not have permissions on.

The controller component of MVC model of Elgg consisting of the *Actions* of the system which are the primary way the users interact with the Elgg site. An action in Elgg Framework is the code that make changes to the database when a user does something such as logging in, posting a comment, and creating an application model. The action script processes input, makes the appropriate modifications to the database, and provides feedback to the user about the action. By default, actions are only available to logged in users and include Cross-Site Request Forgery (CSRF) Security token to overcome session fixation [56], Session Hijacking [57] and Cross-site Scripting [58].

Additional, the controller component includes the *Events* and the *Plugin Hooks*, which are used in Elgg Plugins to interact with the Elgg engine. Events and hooks are triggered at important times throughout Elgg's boot and execution process, and allows plugins to modify or cancel the default behaviour of Elgg. When an event is triggered, a set of handlers is executed in order of priority. Each handler is passed arguments and has a chance to influence the process. After execution, the "trigger" function returns a value based on the behaviour of the handlers.

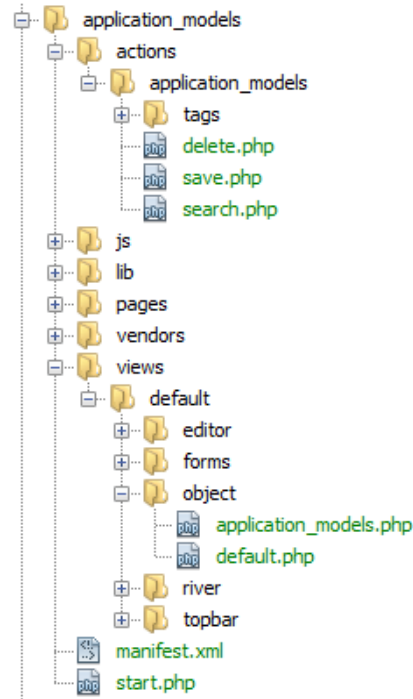
The View component is responsible for creating the output. Generally, this will be HTML sent to a web browser, but it could also be XML, JSON or any other data formats. The Views handles everything from the layout of pages and chunks of presentation output (like a topbar) down to individual links and form inputs.

The extensibility of Elgg can be established not by modifying the core system but by introducing new plug-ins which follow the MVC model. A new plug-in can create a new entity. Thus, each entity is characterized by a numeric Globally Unique Identifier and Access ID. The Access ID determines the permissions that other users have. Thus, when a page requests data, it never touches those data that the current user does not have permission to see. All plug-ins share a common structure of folders and PHP files, following the MVC model of figure 3.2.

The hierarchy of a plug-in is shown in figure 3.4. Folder *actions* includes the actions applied on application models. Every active participation by the user is performed via an action. Logging in, creating, updating or deleting content are all generic categories of actions. The *views* folder contains the *php* forms applied on application models, *river* events (Elgg terminology for live feeds). Views are responsible for creating the output for the client browser. Generally, this will be HTML, but it can be also JSON or other format. *Pages* overrides elements of

core Elgg pages and can be from chunks of presentation output (like sidebars) down to individual html code. The *js* and *lib* folder provides javascript and *php* library functions. Finally, the *vendors* folders include third-party frameworks such as Twitter's bootstrap front-end [59]. The most important file of a plug-in is the *start.php* script, which contains the *page handler*. Page handler is a function manages the plug-in pages enabling custom url redirect to a specific page. The plug-in initialization is also defined in the *start.php* and registers actions, events and determines the views.

Figure 3.4: The structure of the application description plug-in.



New functionality of Elgg Social Network Platform

As described in the previous section, the new functionality of Elgg Social Networking Platform can be introduced by new plugins. The modification of the core system is not a good practice because make the system more difficult to implemented and can not upgrade to the new versions of Elgg framework. Though, the following plugins are implemented:

ApplicationModel. The ApplicationModel plugin has a *page handler* to manage the application Model pages. Also, ApplicationModel has client side JavaScript for manipulating User Interaction and dynamic pages. Furthermore, some php libraries are implemented, for example a library for interaction with CDO client or a library for manipulating Application models.

Components. The Components plugin has a *page handler* to manage the Components pages and the Categories of them and a php library to interact with Chef Supermarket.

CustomView. The CustomView plugin has all the necessary customization of the PaaSage Social Networking Platform (PSNP). All custom views of the system are implemented in this plugin. This plugin overrides all the default views of the Elgg that should be changed and contains client side JavaScript. Furthermore, CustomView has the following seven page handlers: *profile* responsible for profile pages, *avatar* responsible for the photos of user pages, *settings* responsible for the pages of user settings, *friends* responsible for the friends of the user pages, *contact* responsible for the Contact Information of the PSNP, *review* responsible for the reviews of Application Models and *search* responsible for the main search facility of PSNP. Finally, CustomView has all the required *Actions* of the plugin such as the vote up or down, the action of add a review etc.

NotificationSystem. This plugin is responsible for the notifications of the Social Network which contains a relevant page handler, JavaScript for client side scripting and a php server side library.

Tags. This plugin does not include any page handler but only the necessary actions for the Tags such as *add* or *delete* a Tag and a php library responsible for those.

UserStatistics. This plugin is responsible for collecting and displaying the information about the Users.

Memcached. This plugin has all the essential functionality for memcached implementation as will described in section 3.3.

ZookeeperRecipes. This plugin has all the essential functionality for memcached implementation as will described in section 3.2.

Groups. The Groups plugin is the default plugin of Elgg Framework modified to support the required functionality.

Messages. The Messages plugin is the default plugin of Elgg Framework modified to support the required functionality.

Twitter bootstrapping of Elgg

Twitter Bootstrap [60] [61] is a free and open-source collection of tools for creating dynamic websites and web applications. It contains HTML and CSS-based design templates for typography, forms, buttons, navigation and other interface components, as well as optional JavaScript extensions. It aims to ease the development of dynamic websites and web applications.

We customize the Elgg View inserting the Twitter Bootstrap View System. The default View System of Elgg changed to support the Bootstrap responsive grip.

3.1.1 CDO communication with Social Networking Platform

As mentioned in 1.1, the execution history of deployments of application models and the description of those models are stored in the CAMEL information repository, which is implemented as an Eclipse CDO server. In order to communicate with the CDO Server, a CDO Java Client is needed. The CDO client is in the middle between the Social Networking Engine (Elgg Server) and the CDO server, making possible the exchange of information between those two as shown in figure 3.1.

Specifically, for the communication between the CDO Server and the Client, the CDO Client opens one or more sessions to the CDO Server. Each session represents a connection to the CDO repository and provides a broad API to interact with it. A session does not provide direct access to model instances; views or transactions are needed to navigate or modify the model instance graph. The implemented CDO Client expose read/write access to the repository for viewing the execution histories or the model of the applications, or store new execution models. For the communication between the Social Networking Engine and the CDO Client, the CDO Client expose a RESTful API to the Social Networking engine providing all the necessary methods.

For example, when a user from the Social Network requests the execution histories of an application, the Engine sends a request to the CDO Client through the RESTful API, the CDO Client receives the request and forwards an appropriate request to the CDO Server. The CDO Server receives the request and queries the Repository of Application models and executions. When the CDO Server receiving the response, it forwards the response back to the CDO Client, which forwards the response back to the Social Networking Engine. The Social Networking Engine transforms the response to JSON format, in order to be readable by the JavaScript. JavaScript plays the final role, by viewing the execution histories to a proper table to the user that request the page of the executions of application.

3.2 Scaling Social Network Engine

This subsection describes how the horizontal scale of Social Network engine achieved. At the layer 1 of figure 3.2 lives the Apache2 server which as the stress test of the system shows in chapter 4, the Apache2 server takes a heavy load on CPU utilization.

The heavy load of Apache2 server occurs by the nature of Elgg Framework. Because the Elgg core system is implemented to be extensible and configurable, at every time a simple page or just an AJAX call is received by the Elgg, the Elgg Framework performs the following heavy task: broadcasts an *init system* event; this event is caught by all plugins of Elgg and at this initialization phase the plugins registers (1) the page handlers, (2) the php libraries, (3) the actions, (4) the events and hooks, (5) the JavaScript libraries and (6) the CSS scripts. Therefore, we introduce more than one Apache2 Servers running the Social Networking Engine

of Elgg Framework.

The Social Networking Engine keeps some information in the file system instead of the Social Network DB. This information includes the profile photos of users and any other photo such as photos that users add to the community groups. Furthermore, the initial configuration of Social Networking Engine keeps in the file system some caching files representing some views which are independent from specific users and do not change among all users. This file system caching feature removed from the Social Network Engine because is more efficient to use memcached for the caching instead of the slow file system.

The Network File System (NFS) is configured and used in order to all SN Engines have access to the same file system store. An NFS server installed in one of the SN Engines and all the other SN Engines have an NFS client accessing the remote file system.

Distributing Social Network Engine was not an easy problem, so Apache ZooKeeper [51] is used to enable highly reliable distributed coordination among access to the file system storing by Social Networking Engine. Apache ZooKeeper provides a tree abstraction where every node in that tree (or znode) is a file on which a variety of simple operations can be performed. ZooKeeper orders operations on znodes so that they occur atomically. Therefore there is no need to use complex locking protocols to ensure that only one process can access a znode at a time. The tree represents a hierarchical namespace, so that many distinct distributed systems can use a single ZooKeeper instance without worrying about their files having the same name.

Social Networking Engine uses Apache ZooKeeper in order to keep consistent the file system in rear but possible scenarios such as two users try to upload in the same time a file to the same group. When a SN Engine wants to write a file in file system first locks the specific path and after finish the write operation releases the lock.

For communication between the Apache ZooKeeper and the Elgg framework, the `php-zookeeper-recipes` [62] are used by the `ZookeeperRecipes` plugin. Specifically, the *exclusive locks* of Zookeeper are used to keep the system in consistent mode.

3.3 Memcache

This section describes the experience gained by using memcached [63]. Memcached is an open source, high-performance, distributed memory object caching system. We choose memcached, because is a generic simple in-memory key-value store. It has a powerful API available for php. After memcached integration the system increase the response time and the performance.

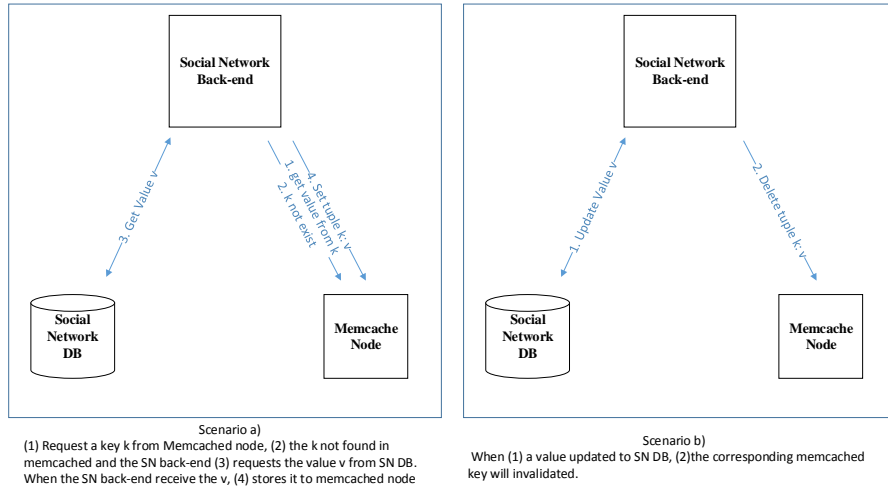
Memcached added in layer 2 of system architecture and used for storing the following key-value tuples: (1) values from Social Network Database such as entities of Social Network, applications, components, users, group discussions, (2) evaluated

javascript code results and (3) executions histories from repository of application models. Storing the executions of applications at Memcached the response time of the system increased because the PHP modules do not need to go through the heavy CDO client but get directly the executions of applications from Memcached.

The apache jmeter [64] was used to measure the response time of the system and the sysstat tool [65] was used to measure the cpu usage. Section 4.1 shows the performance results of this implementation.

All tuples at Memcached are inserted with maximum key expiration time of thirty days. When a value in social network is updated, the memcached key will be deleted as shows the figure 3.5.

Figure 3.5: The scenario a depicts a request from memcached when the key does not exist and scenario b depicts a updated operation of a value)



The Elgg Framework comes with the potentiality to use memcached but is restricted to the facts that a memcached node must to be in the same machine as the Elgg Framework and makes more difficult to configure and strict when an insertion in memcached node will take place. Therefore, a new plugin is implemented called **Memcached** using the memcached php library [66]. The basic memcached functions offered by this plugin are: (1) Add memcached nodes, (2) add a key to a node, (3) get a value of a key and (4) delete a key-value item.

3.4 Natural Language Processing and classification

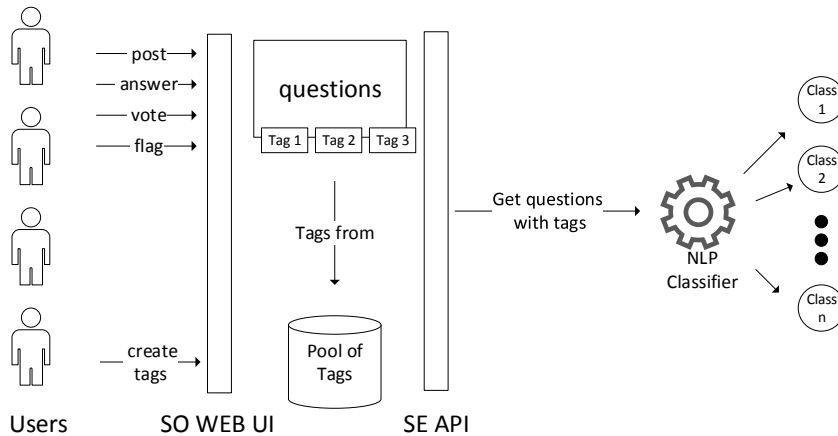
Natural Language Processing(NLP) [36] is a feature added to the Social Network. NLP is used in the interactions between the users and the SN and in the way the platform can understand and determine the type of user input. Particularly, Naïve Bayes Natural Language Understanding algorithm is added to the SN using the

Natural framework [67] implemented with node.js. In general, Machine Learning algorithms such as Naivy Bayes, want an input, calling data set, of training data, this training data is pulled from StackOverflow (SO) questions. Those questions are an excellent repository to train the NB algorithm, because is categorised by tags and in the time that this implementation was took place, the Social Network had not a good repository of many questions.

In general, the main actions of SO members is shown in figure 3.6. When users ask questions in SO, they must specify some tags describing their questions. A tag is a keyword or label that categorizes their question with other, similar questions. The users of SO sometimes, try to add as much tags as possible in order to make their questions popular and get them answered. SO restricts the users to add up to five tags in each of their questions. After a question is posed, the community of SO can vote up or down the question or privileged users can flagged the question as *duplicate*, *off-topic*, *unclear*, *too board* or *primarily opinion-based*. So low quality questions will can be removed from the site and keep the questions repository clear and helpful to other potential developers.

For the training sets, which is the data to train the NB algorithm, the most voted question of the SO community is used. Those questions have emerged as the good questions in their fields and surely, we avoid the case to use a training set with miss-tagged questions which would result a miss-guided NLP classification. The NLP training set is retrieved from SO site using the stack exchange(SE) API [68]. The SE API is a powerful API, which allows to take the questions, the answers, the users and all the information that exists in SO site through a programming interface.

Figure 3.6: The main StackOverflow users' actions and NLP Classifier.



The first training set of our Natural Processing Tool consisted of five tags, relative to our platform. Those tags were: *scalability*, *reliability*, *design*, *performance* and *optimization* with thirty questions per tag. For each of the tags, the thirty

most voted questions from StackOverflow were retrieved and classified to each specific tag. Those exactly tags, after classification, are transformed in classes in NLP classification, as shows the figure 3.6

Every time a user asks a question to a platform's community, the classifier determines the class of the question. Then, if the platform is able to determine a heuristic answer, it will post it to the user's question. All the users of the platform can vote up or down this answer, depending on its accuracy, or provide their own answers.

The second training set of Natural Processing Tool was retrieved by automatically discovering tags. NB is trained with 10K questions from StackOverflow Q&A site, fifty questions per tag. The general algorithm is shown below. Firstly, the algorithm starts with a tag which is relevant to the Social Network Platform such as *scalability* at line 01. Afterwards, using the SE API the algorithm gets the 5 most voted questions tagged with *scalability*. For each question (line 04), the `populateClassifier` clears the body from any html tags inserted by the StackOverflow users to beautify their questions (line 05) and classifies this question's body with each tag. Automatically, the `populateClassifier` proceeds to the next tag of this question. When the `populateClassifier` is finished the NB is trained. It should be noted that a question may have more than one tag, so a question can be classified up to five tags / classes. Changing the threshold parameter at the following algorithm, the *populateClassifier* can classified with an arbitrary number of tags. At the following section the process to Bayes classification is described in more details.

```

01: var tags = [ 'scalability' ]
02: populateClassifier(0)
03: function populateClassifier(index) {
04:   var questions = stackexchange.api.getQuestionsByTag(tags[index])
05:   foreach(questions as q)
06:     body = clear(q.body)
07:     foreach(q.tags as t)
08:       classify(body, t)
09:       if(not tags.exist(t) and tags.length() < threshold)
10:         tags.push(t)
11:         populateClassifier(tags.indexOf(t))
12: }

```

3.4.1 Bayes Classification Algorithm

In this section the Bayes classification algorithm is being described. As the above code snippet shows at line 08, the algorithm classifies a document named *body* into the class *t*. Diving in this function, the *body* is transformed to lower case and the Porter Stemming Algorithm [69] is used for suffix stripping, so the plural part of the words and the suffixes are removed (such as *-ing* and *s*). Thus the following words: *connected*, *connection*, *connections*, *connecting* will all be transformed to the single

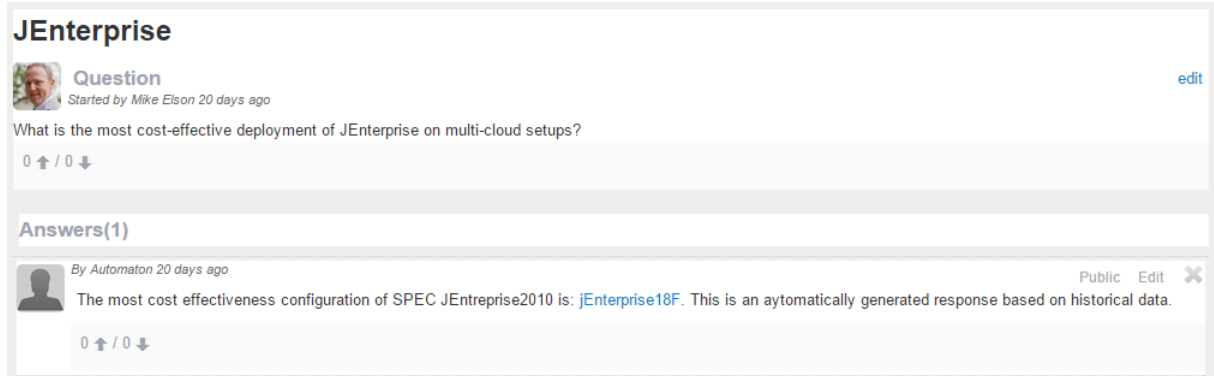


Figure 3.7: Automated answer to user’s question using NLP.

word “connect”. The Porter Stemming Algorithm is not using any dictionary but a simple list of suffixes which makes the algorithm fast (10.000 different words in 8.1 seconds). After this process a table of words of this *body* is kept.

After the `populateClassifier` has finished, the `trainClassifier` is called (for simplicity not shown in the above code snippet). The objective of `trainClassifier` is to make the document body ready for Bayes Classification. So, `trainClassifier` counts the number of occurrences of each word in each class.

Since the Classifier is ready, when a future request for classification comes, the Classifier returns the probability for each class to be part of this class. The probability of each class is calculated with the following formula:

$$prob(d/c) = \log \left(\frac{countedTerms(d, c)}{totalsTerms(c)} \right)$$

Where the probability of a document d to be a class c is the logarithmic value of the division of the words(terms) of d found in class c by the total number of terms in c .

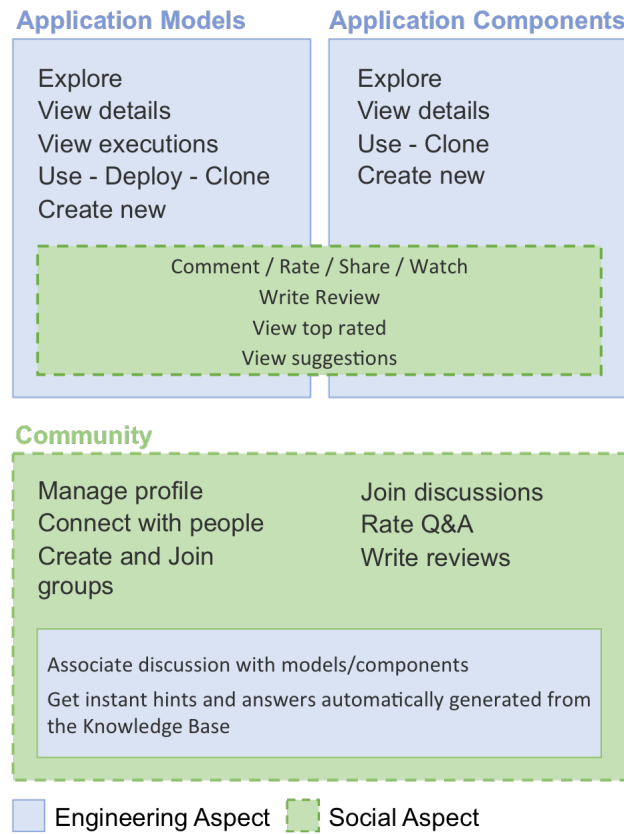
3.4.2 Automated answers from Natural Language Processing

As described in section 3.4 Natural Language Processing is used to determine the users’ input question in groups. When a user asks a question, the body of the question classified into categories and when the NLP classifies the question to a specific category, an approximate answer can be given in response. As shown in figure 3.7, after the user’s question, the classifier process the body and if the question is about the *JEnterprise* and the *cost effectiveness* the approxiame answer “The most cost effectiveness configuration of SPEC JEnterprise2010 is: `jEnterprise18F` ...” is given. The users of the PaaSage Social Network Platform can vote up or down the answer and/or provide their own answers.

3.5 User interface

In this section described the User Interface of Social Network implemented based on 104 mock-ups created by HCI expert team. In order to support those look & feel and the functionality of those mock-ups 25K lines of php, js and css code is written. The key design objective of the social network platform is to create a strong bond between (i) software engineering services for managing and deploying cloud-targeted application models; and (ii) community-oriented facilities for communication and collaboration between users. The interconnections between the two in the design of the user interface are depicted in Figure 3.8. The prototype implementation is publicly accessible on-line at <http://socialnetwork.paasage.eu>.

Figure 3.8: The engineering & social activities are seamlessly within the Platform.



3.5.1 User Interface Design Principles

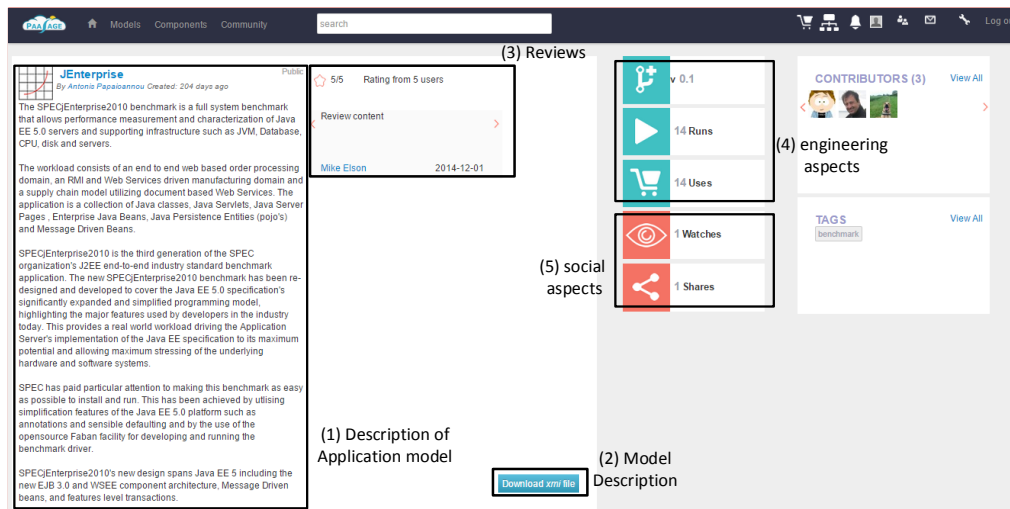
The discrete entities, which bound together the Social Networking with application model aspects of Platform are:

- *Application Models.* Application Models is a key entity of the platform. An

example is shown in figure 3.9, consisting of a human friendly description (label 1 in fig.3.9), the Camel Description of the model (label 2 in fig.3.9), reviews about the model (label 3 in fig.3.9). An overview of engineering aspects such as version and runs (label 4 in fig.3.9) and an overview of social aspects such as share and watch (label 5 in fig.3.9). The *share* action broadcast the model to the friends of the user that shares the model. The *watch* action notifies the user for future updates of the application model.

- *Components.* We have integrated the Chef supermarket components into Social Network Platform. The components help the DevOps users to generate their application models as described in 3.5.2.
- *Users.* Users which basically be Cloud Deployment specialists and other users who want to know which deployment configuration should use. Users can exchange knowledge to groups and benefit from CAMEL repository. They can create or join groups, ask and answer questions, follow application models and create their own network of friends.
- *Groups.* Every user of PaaSage Social Network Platform can create or join groups. Groups help users to interact with each other and gain knowledge from experts.

Figure 3.9: The application model home page



Gamification

Following recent trends in social networks design and with the aim to motivate users active and regular participation in the professional network, the design em-

plays gamification features, namely use of video game elements to improve user experience and user engagement in non-game services and applications [70]. One gamification feature in the Social Network design is the reward system for active community members. As users contribute content (models, components, ratings, reviews, questions, or answers) they receive experience points leading to special badges visible to all community members. Other features are the Profile completeness bar with suggestions on how to increase it. Finally, the concept of Model badges awarded to application and component models in case of excelling performance. Badges can serve among others as goal-setting devices, status symbols, and indications of reputation assessment procedures [71].

3.5.2 Automated Application Model Creation

The DevOps users of PaaSage Social Network Platform can benefit from the automated creation of CAMEL baseline models or upload their own created models using external editors like EMF [10] tree based editor or GMF [72] editor.

In order users to create automated generated baseline CAMEL models, they can browse around the integrated Chef Components inside the Social Network and find the appropriate components for their applications. The platform has integrate a pointer for each Chef cookbook to the Social Network Database using the Chef Supermarket API [73]. A PHP command has implemented in order to iterate through all Chef cookbooks and update the repository of SNP. This command has been configured and runs one time every day.

Through the Application Model Creation Page of the PaaSage Social Network Platform, the users can upload an external CAMEL model description of their application or create a new one with the help of the Platform in 4 simple steps as shown in figure 3.10. In the step zero 3.10a, the user asked if the new Application model has already an CAMEL model or the user want to create a new Application CAMEL model through automated generation. For automated generation the users should have put in the user's cart the Chef cookbooks that wants. Then, in step one 3.10b, the user selects from his/her list of components which of those will be included in the Application Model. As the figure 3.10b shows, the user has four components *mysqld*, *apache2*, *nodejs* and *ruby_installer* and selects three of them. In the next step, shown in figure 3.10c the user provides the deployment information (to which cloud provider the Application will run and which type of VMs will be used). Also, some components, in this example the *nodejs* component will be collocated in the same VM as the *apache2*. In step tree, shown in figure 3.10d, the user provides the communication information between the components, for example the *nodejs* communicates with *mysqld* in the default mysql port *3306*.

In the final step as shown in figure 3.11, the user provides the final needed information about the name of the Application model, a human friendly description and the version of the model. In the bottom of the form of the figure 3.11 the user can find three actions: the *Previous* action, the *Save as draft* action and the *Finalize* action. The *Previous* action can performed all around the steps and make

the user easily walk around the steps. The *Save as draft* action creates the CAMEL model of the application but the model is publicly available and the creator or the contributors of the Application can edit the Model. At the *Finalize* action, the model is now publicly available to the PaaSage Social Network users.

New Application Model

Step 0 Step 1 Step 2 Step 3 Step 4

A new application model can be:

- A CAMEL abstract model that can be used as input to the PaaSage platform
- A baseline CAMEL concrete (deployable) model created automatically from a list of components.

Currently you have 4 component(s) in your cart.

Note: You can extend the baseline model using EMF editor

Import abstract model Create baseline model

(a) Step 0: Upload external or create baseline model.

New Application Model

Step 0 Step 1 Step 2 Step 3 Step 4

Which of the following components will be associated with the application model?

☒ mysql
☒ apache2
☒ nodejs
☐ ruby_installer

Previous Next

(b) Step 1: Choose the Components from the users' list.

New Application Model

Step 0 Step 1 Step 2 Step 3 Step 4

Select where your components will be deployed

Component	Deployed on		
mysql	VM	Amazon	m1.small
apache2	VM	Amazon	m1.small
nodejs	Colocated	apache2	

Previous Next

(c) Step 2: Deployment information.

New Application Model

Step 0 Step 1 Step 2 Step 3 Step 4

Describe communication connections between components.

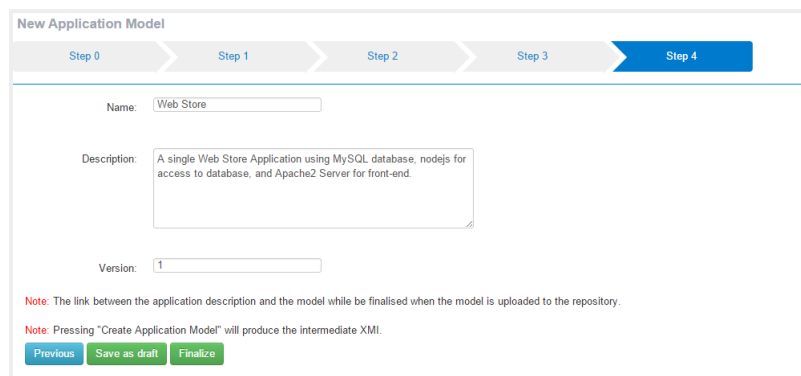
Component	Component	port	mandatory
mysql	nodejs	3306	<input checked="" type="checkbox"/>
nodejs	apache2	3000	<input checked="" type="checkbox"/>

+ Add Connection

Previous Next

(d) Step 3: Communication Information.

Figure 3.10: Steps for Automated creation of baseline model.



New Application Model

Step 0 Step 1 Step 2 Step 3 Step 4

Name:

Description:

Version:

Note: The link between the application description and the model will be finalised when the model is uploaded to the repository.

Note: Pressing "Create Application Model" will produce the intermediate XML.

[Previous](#) [Save as draft](#) [Finalize](#)

Figure 3.11: Final Step of Automated creation of baseline model.

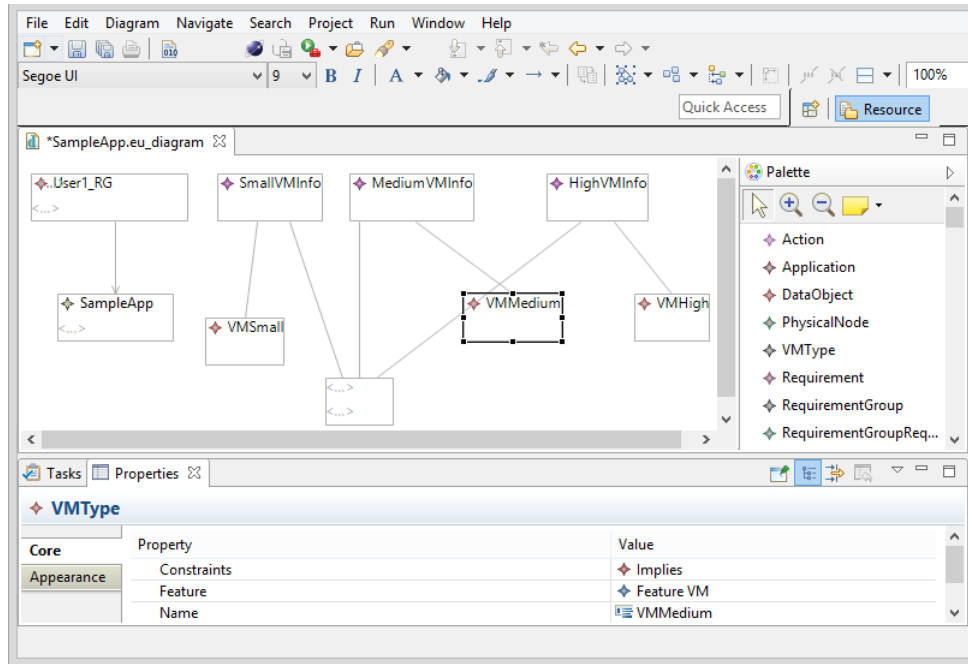


Figure 3.12: GMF editor composition of a sample application.

3.5.3 Graphical modeling of applications

Advanced users of the PaaSage Social Network Platform can compose application models through Graphical Modeling Framework (GMF) [72] which is an external Eclipse editor. GMF provides a set of generative components and runtime infrastructures for developing graphical editors based on Eclipse Modeling Framework (EMF) and Graphical Editing Framework (GEF). The GMF editor is generated from CAMEL *ecore* schema and provides the graphical palette to compose applications.

Figure 3.12 shows the composition of a sample application model with the GMF editor. The palette in the right, contains all nodes and relationships needed to describe an application model. In the center, the composition of a sample application is shown, consisting of three VM types, the VM information about these VMs and the owner/user of the application. The GMF editor generates two files, one responsible for the graphical representation and the XMI description of the application model that can be uploaded to the social networking platform.

Chapter 4

Evaluation

This chapter describes the evaluation of the two different implementations of the system architecture. In the first architecture, more than one memcached instances at layer 2 were introduced, as figure 3.1 shows. In the second, more than one Social Network engines were introduced at layer 1.

In order to measure the response time (RT) the Apache JMeter application [64] is used. The Apache JMeter is an open source benchmark designed to test functional behaviour and measure performance, targeting web applications. Notably, the RT measured by JMeter may not be the real one, because the JMeter measures the elapsed time from just before sending the request to just after the last response from the server has been received. As a result, the time to render the web page to the client web browser and the execution time of JavaScript code is not measured. Because those two time intervals are client limited and depend on client performance and on which web browser is used, they are excluded from the following performance test benches. For the next experiments, a specific web page will be used. This page does not use any AJAX call, in order to not misguide the results. Therefore, the RT measures the time from just before JMeter sends the request to just after the last response is received. During this measured time interval, the Social Network engine performs the following actions:

- I The Social Network engine sends a request to CDO Client for the application execution model.
- II The CDO Client forwards this request to CDO Server.
- III Afterwards CDO Server queries the mysql repository of application models and executions, and finally gets the executions results.
- IV CDO Server forwards the results through the CDO Client to the Social Network Engine.
- V Finally, the Social Network engine sends queries to the Social Network DB in order to get all the necessary Social features for this application page

4.1 Improving Performance with memcached

By adding a memcached node at the system architecture, the Social Network Engine first asks the memcached node if it has the tuples that the SN Engine needs. So the steps(*I* to *V*), mentioned previously, are not necessary if the memcached node has cached the values that the Social Network Engine needs. The loop through CDO Client - CDO Server and the repositories is bypassed.

For the following experiments all the memcached nodes are warm up and have cached all the needed CDO and Social entities information. Furthermore the CDO server has been warm up after a fresh restart. As the table 4.1 shows the starting process of the CDO server produces 1938 queries to MySQL database. The *fresh query* for an application model (both social information and executions) produces 15182 queries to MySQL database. The CDO server caches the results, so a second query for this application model produces 251 queries, which the most of them are the queries for the social information of the application. Introducing memcached, if the request for the application model is cached, the queries to database are lowering to 147.

The test performed with the following loads: (L1) ten users requests *two* applications, (L2) ten users requests *four* applications and (L3) ten users requests *eight* applications. All three Loads run consecutively one hundred times each. Those Loads request applications, which have ten execution rows pulled from the repository of applications models and executions, and about one hundred queries to the Social Network DB. In this experiment we kept constant the following components of the system: the Elgg front-end Apache2 server, the Social network BD, and the CDO server - client communication but increased the number of memcached nodes. The figure 4.1 shows the average, minimum and maximum response time (RT) in milliseconds with the following system configuration: (C1) no memcached node, (C2) one memcached node and (C3) two memcached nodes.

As we going from C1 to C3 and specifically for L(oad) 3, the RT is reduced by 80,4% at C2 and by 88,78% at C3. As the figure 4.1 shows, at the first configuration C1, the L3 takes 8836 ms, an RT which is definitely prohibitive for web applications. Introducing more memcached nodes at C2 and C3 the RT is decreased dramatically at 1764 ms at L2 and at 992 ms at L3. Going from C1 to C2, the 80,4% reduction of RT is due to the introduction of memcached node and bypassed the steps I - V. Going from C2 to C3, the 43,77% reduction of RT is due to adding more memcached nodes, resulting to more cpu cores introduced to the architecture.

Furthermore, the CPU utilization is measured using the sysstat tool [65]. We measured the CPU utilization for all the VMs running the experiment. The information about the VM resources is listed in the table ???. The Social Network Engine and the CDO Client were running at t1.micro instance. The mysql (repositories) and the CDO Server were running at m1.xlarge. The average CPU utilization is shown in the figure 4.2. At the simple configuration C1, even in small loads such as L1, the SN Engine reached 50,39% CPU utilization. In the medium load L2 and big load L3 the SN Engine is kneeled down to 91,96% and 92,86%. This big

Figure 4.1: The average response time for all configurations.

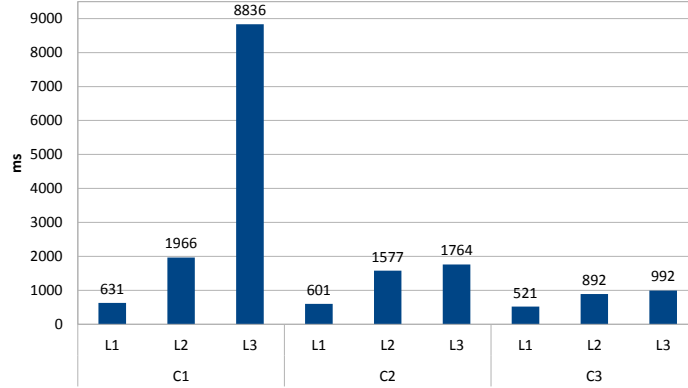


Table 4.1: Number of Queries to Social Network and CDO server Databases.

State	# of Queries
Fresh start	1938
Fresh Query	15182
Cached Query from CDO	251
Cached Query from Memcached	147

consumption of CPU was due to all the initialization that Elgg Social Network Engine has to do for each request and due to CDO Server queries.

Moving from configuration C1 to C2, the CPU consumption went to memcached node. Thus, the Social Network engine was de-congested and the RT improved. However, for the big load L3 the memcached node reached 89,79%. To solve memcached CPU overhead, one more memcached node was added at configuration C3. This second memcached node shared the CPU overhead with the first memcached node and the RT improved furthermore. For all three loads at C3, the first memcached node has more CPU utilization from the second by an approximately factor of 2,2. This difference between the two memcached nodes appeared due to the first node storing more popular key-value pairs than the other.

4.2 Improving Performance with engine

This section evaluates the horizontal scale of Social Network engine as described in 3.2. A memcached node was living between the Social Networking Engines

Figure 4.2: The average CPU utilization for all components.

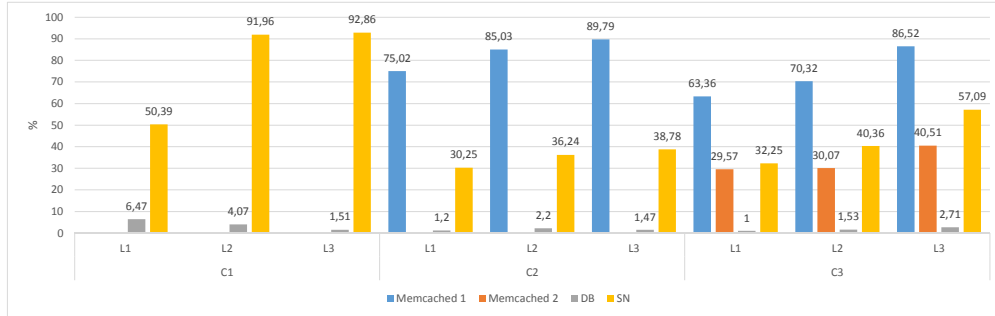


Table 4.2: VM resources

Component	VM type
SN engine, CDO client	t1.micro
memcached	t1.micro
repositories, CDO Server	m1.xlarge
jmeter	m1.large

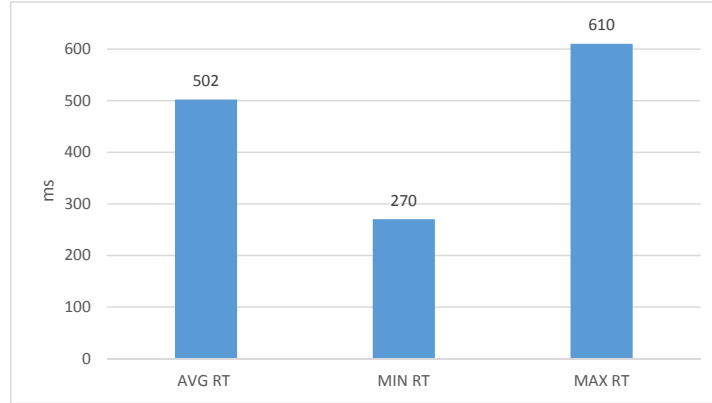
and the back end system. The VM resources are kept the same in the previous experiment, shown in the table ?? . One more Social Networking Engine instance was added with the same type as former SN Engine. The load from one SN Engine is now distributed to two SN Engines instances. So, the response time improved, as shown in the figure 4.3 for the Load 3 compared to the previous test-bench.

The CPU utilization to SN Engines decreased as shown in the figure 4.4. This reduction is due to the requests being distributed to two instances instead of only one. The CPU utilization of the memcached node increased, but this can be solved by introducing more memcached nodes as the previous section describes.

4.3 Evaluation of NLP classification

This section evaluates the Natural Language Processing Tool as described in section 3.4 for the first training set. This set had five different classes or, according to StackOverflow dialect, five tags, as described in 3.4 and shown in table 4.3. For the testing set, thirty questions from StackOverflow were used per class. Those questions were different from the training set and had the highest activity dur-

Figure 4.3: The Response time for two Social Network Engines.



ing that specific time period. Each row at table 4.3 shows a class as classified by StackOverflow users, and each column shows how our classifier classifies the question. For example, twenty one questions about *reliability* were classified correctly, but eight were wrongly classified as *optimization* and one was wrongly classified as *performance*. The misclassification was not an error of our classifier but some questions from the testing set were either wrongly classified by the StackOverflow users or had more than one tags. For example, one question that was wrongly classified as *performance* instead of *reliability* was: “can somebody explain me how to handle errors. my code is: ...”. The above question is out of the scope of reliability and even though the user tagged it as a reliability question, it was downvoted and marked as “very low quality” from the StackOverflow community.

Another example showing that the classifier is not misclassifying is the following question: “I want to perform some data manipulation tasks and analysis in spark and want to **optimize** the run times. here is the problem: ...”. This question was marked by the user with both scalability and optimization tags. Even though this question is retrieved with the scalability tag and in the table it is shown as a wrong classification, after examining this question, one can realize that the *scalability* tag was wrongly added by the user and only the *optimization* tag should be placed.

This shows that our tool can further be used by StackOverflow to mark new questions that are wrongly tagged or misguided. There is a trend among StackOverflow users to add as many tags as they can in order to attract the attention of other users, increase the views of their question and finally get their answers.

Figure 4.4: The CPU utilization for two Social Network Engines.

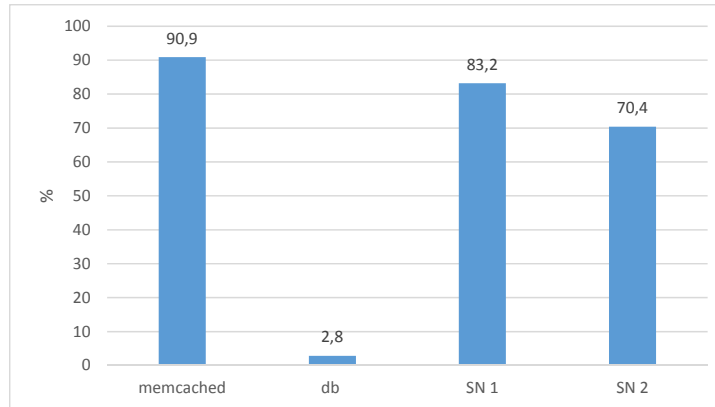


Table 4.3: NLP Evaluation of Classification

class / class	reliability	design	optimazation	performance	scalability
reliability	21	0	8	1	0
design	2	15	8	3	2
optimazation	2	3	15	9	1
performance	2	1	17	9	1
scalability	10	6	3	3	8

4.4 Evaluation of requirements and UI interface

The User Interface of SNP is

Chapter 5

Conclusions and Future Work

In this thesis, a scalable Social Network for DevOps engineers is presented. The scaling of the system presented at the front end layer by introducing more than one Social Networking Engine instances, and/or at the back end of the system by introducing more than one memcached nodes. The DevOps users of the Social Network can benefit from the community knowledge and from the CAMEL repository of application models and executions, to improve the configuration, deployment, and optimization of distributed multi-cloud applications, tasks of major interest to cloud deployment specialists. The design of our professional network applied best practices aiming to support the creation of a vigorous community, to allow users to retrieve timely and appropriate information and to carry out actions in a small number of steps.

Bibliography

- [1] M. Loukides, *What is DevOps?* " O'Reilly Media, Inc.", 2012.
- [2] J. Rossberg, "Collaboration," in *Beginning Application Lifecycle Management*. Springer, 2014, pp. 135–143.
- [3] R. Buyya, M. Pathan, and A. Vakali, *Content delivery networks*. Springer Science & Business Media, 2008, vol. 9.
- [4] P. D. 2.1.2, "Model Based Cloud Platform Upperware," http://www.paasage.eu/images/documents/paasage_d2.1.2_final.pdf, 2014.
- [5] N. Ferry, A. Rossini, F. Chauvel, B. Morin, and A. Solberg, "Towards model-driven provisioning, deployment, monitoring, and adaptation of multi-cloud systems," in *Proceedings of CLOUD 2013: 6th IEEE International Conference on Cloud Computing*, L. O'Conner, Ed. IEEE Computer Society, 2013, pp. 887–894.
- [6] F. Chauvel, N. Ferry, B. Morin, A. Rossini, and A. Solberg, "Models@ runtime to support the iterative and continuous design of autonomic reasoners." in *MODELS@ Run. time*, 2013, pp. 26–38.
- [7] A. Andrieux, K. Czajkowski, A. Dan, K. Keahey, H. Ludwig, T. Nakata, J. Pruyne, J. Rofrano, S. Tuecke, and M. Xu, "Web services agreement specification (ws-agreement)," in *Open Grid Forum*, vol. 128, 2007, p. 216.
- [8] C. Quinton, N. Haderer, R. Rouvoy, and L. Duchien, "Towards multi-cloud configurations using feature models and ontologies," in *Proceedings of the 2013 international workshop on Multi-cloud applications and federated clouds*. ACM, 2013, pp. 21–26.
- [9] D. Steinberg, F. Budinsky, E. Merks, and M. Paternostro, *EMF: eclipse modeling framework*. Pearson Education, 2008.
- [10] Eclipse, "CDO Model Repository," <http://projects.eclipse.org/projects/modeling.emf.cdo>, 2015, [Online; accessed 29-July-2015].
- [11] K. Kritikos, M. Korozi, B. Kryza, T. Kirkham, A. Leonidis, K. Magoutis, P. Massonet, S. Ntoa, A. Papaioannou, C. Papoulas, C. Sheridan, and

- C. Zeginis, “D4.1.1 – prototype metadata database and social network,” Accessed 8/2015, available from http://www.paasage.eu/images/documents/PaaSage-D4.1.1_final.pdf.
- [12] A. Papaioannou and K. Magoutis, “An Architecture for Evaluating Distributed Application Deployments in Multi-Clouds,” in *Proceedings of 5th IEEE International Conference on Cloud Computing Technology and Science (CloudCom’13)*. Bristol, UK: IEEE, 2013.
 - [13] PaaSage EU FP7 project, <http://www.paasage.eu/>, [Online; accessed 29-July-2015].
 - [14] D. Baur, S. Wesner, and J. Domaschka, “Towards a model-based executionware for deploying multi-cloud applications,” in *Advances in Service-Oriented and Cloud Computing*. Springer, 2014, pp. 124–138.
 - [15] R. Sumbaly, J. Kreps, L. Gao, A. Feinberg, C. Soman, and S. Shah, “Serving large-scale batch computed data with project voldemort,” in *Proceedings of the 10th USENIX conference on File and Storage Technologies*. USENIX Association, 2012, pp. 18–18.
 - [16] G. DeCandia, D. Hastorun, M. Jampani, G. Kakulapati, A. Lakshman, A. Pilchin, S. Sivasubramanian, P. Voshall, and W. Vogels, “Dynamo: amazon’s highly available key-value store,” in *ACM SIGOPS Operating Systems Review*, vol. 41, no. 6. ACM, 2007, pp. 205–220.
 - [17] R. Nishtala, H. Fugal, S. Grimm, M. Kwiatkowski, H. Lee, H. C. Li, R. McElroy, M. Paleczny, D. Peek, P. Saab *et al.*, “Scaling memcache at facebook.” in *nsdi*, vol. 13, 2013, pp. 385–398.
 - [18] D. W. Jorgenson and K. J. Stiroh, “Information technology and growth,” *American Economic Review*, pp. 109–115, 1999.
 - [19] “Social networks popular among programmers,” Accessed 8/2015. [Online]. Available: <http://www.informationweek.com/wireless/social-networks-popular-among-programmers/d/d-id/1078472>
 - [20] “Github,” <http://github.com/>, [Online; accessed 24-May-2015].
 - [21] “Gitter: The chat for github,” <http://gitter.im/>, [Online; accessed 24-May-2015].
 - [22] “Sourceforge,” <http://sourceforge.net>, [Online; accessed 24-May-2015].
 - [23] “Google code,” <https://code.google.com>, [Online; accessed 24-May-2015].
 - [24] “Codeplex: Project hosting for open source software,” <https://www.codeplex.com>, [Online; accessed 24-May-2015].

- [25] “Stackoverflow,” <http://stackoverflow.com/>, [Online; accessed 24-May-2015].
- [26] B. Vasilescu, V. Filkov, and A. Serebrenik, “Stackoverflow and github: Associations between software development and crowdsourced knowledge,” in *Social Computing (SocialCom), 2013 International Conference on*, Sept 2013, pp. 188–195.
- [27] “Ibm blue mix for developers,” <http://https://developer.ibm.com/bluemix/>, [Online; accessed 24-May-2015].
- [28] “Ibm devops best practices,” <http://www.ibm.com/developerworks/devops/practices.html>, [Online; accessed 24-May-2015].
- [29] R. Sumbaly, J. Kreps, and S. Shah, “The big data ecosystem at linkedin,” in *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*. ACM, 2013, pp. 1125–1134.
- [30] “Geeklist,” <https://geekli.st/>, [Online; accessed 10-June-2015].
- [31] “Snipplr,” <http://snipplr.com/>, [Online; accessed 10-June-2015].
- [32] “Masterbranch,” <https://masterbranch.com/>, [Online; accessed 10-June-2015].
- [33] “Dzone,” <http://www.dzone.com/>, [Online; accessed 10-June-2015].
- [34] “The code project,” <http://www.codeproject.com/>, [Online; accessed 10-June-2015].
- [35] J. Howe, “The rise of crowdsourcing,” *Wired magazine*, vol. 14, no. 6, pp. 1–4, 2006.
- [36] C. D. Manning and H. Schütze, *Foundations of statistical natural language processing*. MIT press, 1999.
- [37] A. Pak and P. Paroubek, “Twitter as a corpus for sentiment analysis and opinion mining,” in *LREC*, vol. 10, 2010, pp. 1320–1326.
- [38] S. Verma, S. Vieweg, W. J. Corvey, L. Palen, J. H. Martin, M. Palmer, A. Schram, and K. M. Anderson, “Natural language processing to the rescue? extracting" situational awareness" tweets during mass emergency.” in *ICWSM*. Citeseer, 2011.
- [39] A. Go, R. Bhayani, and L. Huang, “Twitter sentiment classification using distant supervision,” *CS224N Project Report, Stanford*, vol. 1, p. 12, 2009.
- [40] E. Wong, J. Yang, and L. Tan, “Autocomment: Mining question and answer sites for automatic comment generation,” in *Automated Software Engineering (ASE), 2013 IEEE/ACM 28th International Conference on*. IEEE, 2013, pp. 562–567.

- [41] C. Lueninghoener, “Getting started with configuration management,” 2011.
- [42] Bcfg2, Accessed 8/2015, <http://www.bcfg2.org/>.
- [43] CFEngine, Accessed 8/2015, <http://www.cfengine.com/>.
- [44] Chef, Accessed 8/2015, <https://www.chef.io/>.
- [45] Puppet, Accessed 8/2015, <http://www.puppetlabs.com/>.
- [46] A. Tsalolikhin, “Summary, configuration management summit,” 2010.
- [47] T. Delaet, W. Joosen, and B. Vanbrabant, “A survey of system configuration tools,” in *Proceedings of the 24th International Conference on Large Installation System Administration (LISA ’10)*. San Jose, CA: ACM, 11/2010, pp. 1–8.
- [48] M. Fowler and M. Foemmel, “Continuous integration,” *Thought-Works*) <http://www.thoughtworks.com/continuous-integration>, 2006.
- [49] “Travis continuous integration,” Accessed 8/2015, <https://travis-ci.org/>.
- [50] “Jenkins continuous integration,” Accessed 8/2015, <https://jenkins-ci.org/>.
- [51] A. Zookeeper, “Apache zookeeper,” <https://zookeeper.apache.org/>, 2015, [Online; accessed 20-May-2015].
- [52] E. S. N. Engine, “Elgg social networking engine,” <http://elgg.org/>, 2015, [Online; accessed 19-May-2015].
- [53] jQuery, “jquery,” <https://jquery.com/>, 2015, [Online; accessed 19-May-2015].
- [54] E. McCormick and K. De Volder, “Jquery: finding your way through tangled code,” in *Companion to the 19th annual ACM SIGPLAN conference on Object-oriented programming systems, languages, and applications*. ACM, 2004, pp. 9–10.
- [55] “Chart.js, simple, clean and engaging charts for designers and developers,” Accessed 9/2015, <http://www.chartjs.org/>.
- [56] M. Kolšek, “Session fixation vulnerability in web-based applications,” *Acros Security*, p. 7, 2002.
- [57] W. Burgers, R. Verdult, and M. van Eekelen, “Poster: Prevent session hijacking.”
- [58] J. L. Thames, “Comparing cross-site scripting vulnerabilities.”
- [59] “Bootstrap front-end framework,” <http://getbootstrap.com/2.3.2/>, [Online; accessed 24-May-2015].

- [60] “Twitter Bootstrap,” <http://getbootstrap.com/>, 2015, [Online; accessed 12-Aug-2015].
- [61] D. Cochran, *Twitter Bootstrap Web Development How-To*. Packt Publishing Ltd, 2012.
- [62] “The Zookeeper recipes,” <https://github.com/Gutza/php-zookeeper-recipes>, 2015, [Online; accessed 12-Aug-2015].
- [63] Memcache, “Memcache,” <http://memcached.org/>, 2015, [Online; accessed 18-May-2015].
- [64] A. jMeter, “Apache jmeter,” <http://jmeter.apache.org/>, 2015, [Online; accessed 19-May-2015].
- [65] sysstat, “Performance monitoring tools for linux,” <https://github.com/sysstat/sysstat>, 2015, [Online; accessed 19-May-2015].
- [66] “PHP Memcached library documentation,” <http://php.net/manual/en/book.memcached.php>, 2015, [Online; accessed 12-Aug-2015].
- [67] “Natural language processing with node.js,” <https://github.com/NaturalNode/natural/>, [Online; accessed 1-July-2015].
- [68] “Stack exchange application programming interface,” <https://api.stackexchange.com/>, [Online; accessed 14-July-2015].
- [69] M. F. Porter, “An algorithm for suffix stripping,” *Program*, vol. 14, no. 3, pp. 130–137, 1980.
- [70] S. Deterding, M. Sicart, L. Nacke, K. O’Hara, and D. Dixon, “Gamification. using game-design elements in non-gaming contexts,” in *CHI’11 Extended Abstracts on Human Factors in Computing Systems*. ACM, 2011, pp. 2425–2428.
- [71] J. Antin and E. F. Churchill, “Badges in social media: A social psychological perspective,” in *CHI 2011 Gamification Workshop Proceedings (Vancouver, BC, Canada, 2011)*, 2011.
- [72] “Graphical Modeling Framework,” https://wiki.eclipse.org/Graphical_Modeling_Framework, 2015, [Online; accessed 30-July-2015].
- [73] “Cookbooks Site API,” https://docs.chef.io/api_cookbooks_site.html, 2015, [Online; accessed 30-July-2015].