

University of Warsaw
Faculty of Mathematics, Informatics and Mechanics

Piotr Styczyński

Student no. 386038

Michał Balcerzak

Student no. 385130

Michał Ołtarzewski

Student no. 382783

Gor Safaryan

Student no. 381501

AWS Cost Optimization Tool

**Bachelor's thesis
in COMPUTER SCIENCE**

Supervisor:
dr Janina Mincer-Daszkiewicz
Instytut Informatyki

February 2019

Supervisor's statement

Hereby I confirm that the presented thesis was prepared under my supervision and that it fulfils the requirements for the degree of Bachelor of Computer Science.

Date

Supervisor's signature

Authors' statements

Hereby I declare that the presented thesis was prepared by me and none of its contents was obtained by means that are against the law.

The thesis has never before been a subject of any procedure of obtaining an academic degree.

Moreover, I declare that the present version of the thesis is identical to the attached electronic version.

Date

Authors' signatures

Abstract

Authors describe the design and implementation process, in-depth system and code architecture as well as used communication protocols, storing methods and other internals of the AWS Cost Optimization System.

Keywords

AWS, Amazon Web Services, cloud computing, cost optimization, cost management

Thesis domain (Socrates-Erasmus subject area codes)

11.3 Informatics, Computer Science

Subject classification

D. Software

Tytuł pracy w języku polskim

Narzędzie do Optymalizacji Kosztów AWS

Contents

1. Introduction
 - 1.1 Overview
 - 1.2 Aim of the thesis
 - 1.3 Structure of the thesis
 - 1.4 Contribution of each author
2. Problem statement
 - 2.1 Motivation
 - 2.2 Overview of existing solutions
 - 2.3 Our solution
3. Project development
 - 3.1 Version control
 - 3.2 Continuous integration
 - 3.3 Communication
 - 3.4 Workflow
4. Tool for AWS cost optimization
 - 4.1 Overall architecture
 - 4.2 Technology stack
 - 4.3 Data models
 - 4.4 Views and design
 - 4.5 API
 - 4.6 Communication integration
 - 4.7 Configuration
5. Summary
- A Deployment and integration guide

Chapter 1

Introduction

1.1. Overview

Cloud computing has become recently one of the most important paradigm shifts in the area of real world software engineering. It has reshaped the whole process of how applications are developed and reduced the amount of upfront investment required to start an internet business. While commercial cloud computing services were first offered in 2006 by Amazon Inc, the original idea and preliminary implementation traces back to Multics OS developed by MIT, GE and Bell Labs. However the idea of time-sharing systems that was the ancestor of further cloud concept was widespread in 60ies [Markus].

The term “Cloud computing” can refer to every layer of application stack: hardware, hosting platform, software and even to a single function. Cloud computing refers to both the applications delivered as services over the Internet and the hardware and systems software in the data centers that provide these services. The services themselves have long been referred to as Software as a Service (SaaS) [Armbrust]. Generally speaking cloud can be perceived as a shared pool of computer resources such as computing capacity, transient and persistent memory, which can be acquired or released on demand. The undisputed power of cloud computing constitutes in its elasticity and granularity: it allows users to ask for hundreds of computers for only 5 minute usage which are shipped during several minutes. Such services are usually offered over remote network connection and users are billed for the portion of the resources they have used. Depending on the cloud infrastructure type the payment models can be different [Laatikainen], but the common spendings are associated with data storage, data transfer, and computing timeshare.

Nowadays the industry increasingly relies on cloud technologies. More and more companies start or move their products to cloud environment. Unfortunately it comes with additional financial costs imposed by cloud providers. In order to make business profitable companies try to reduce amount of money they are supposed to pay to the minimum. Despite different solutions, like employing specific cloud cost optimization team, more and more firms decide to benefit from dedicated software, which is supposed to help manage and optimize their cloud usage. It is not easy to choose the right tool having that many choices.

There are some notable solutions of AWS cloud optimization problem – *AWS Cost Explorer* and *AWS Cost Management*, *Cloudability*, *Apptio* and others widely used nowadays. In spite of that fact there is still place for new tools targeting omitted types of clients or wrapping and bundling the greatest features from existing ones.

1.2. Aim of the thesis

The primary objective of the thesis is to create a tool complementing existing solutions used in cloud cost optimisation. We focus on Amazon Web Services platform maintained by Amazon as it is one of the most commonly used. In our tool, called *Omigost*, we will try to target small/middle sized companies creating simple and easy to use software with flexible configuration options.

1.3. Structure of the thesis

The thesis is structured as follows. In Chapter 2 we describe the problem of cloud cost optimization with additional description of selected existing solutions. Then, in Chapter 3 we present our solution in detail. We mention, among others, whole system architecture, API and configuration. Finally, in Chapter 4 we sum up the whole thesis. In Appendix A we describe how to introduce our solution in a company.

1.4. Contribution of each author

It is important to mention that each author worked to some extent on every part of the thesis. However authors contributed mainly to the following parts:

- Michał Ołtarzewski
 - Software development process management
 - Design of backend architecture
 - Backend part implementation
- Michał Balcerzak
 - Research of existing solutions
 - Design of backend architecture
 - Backend part implementation
- Piotr Styczyński
 - Design and prototype of visual part of tool
 - Frontend part implementation
- Gor Safaryan
 - Research of AWS APIs and SDKs
 - Design of backend architecture
 - Backend part implementation

Chapter 2

Problem statement

2.1. Motivation

As there are plenty of various billing models for cloud services [**GLaatikainen**], the effective management of them became a tough problem. The ease of resource allocation led to situation when tracking tiniest details of billings is an unaffordable challenge.

The tooling that exists is targetting wideworld-scale companies that are able to require expensive licenses and hire cost-optimization teams. The software as it is in case of Cloudability is too complex for average user and does not provide easy way to incorporate custom business flows into the tool. Amazon, as one of the leading cloud providers, offers different tools for exploration of expenses. Most commonly used options are either their public APIs [**AWSCostManagement**] or specific SDKs supporting lots of languages. Unfortunately previously mentioned tools are rather simple and do not satisfy all the needs of potential clients.

The common case that is unresolved is the distribution of research and development resources. We observed that there exists no tool that would support request for resources of individual worker with regards to custom management propagations as specified by client bussiness model. Cloudability [**CloudabilityAlerts**] offers simple alerts, but they lack Slack support and beforementioned propagation abilities.

There exists an obvious gap in the market, that our solution, *Omigost*, will try to cover. Having one versatile tool removes the need for using few detached pieces of software. The most important aspects of our tool are intuitive interface and different types of user notifications, which highly increases spendings clarity and helps in decisions connected to costs cutting. We hope it will allow companies to focus more on providing value to their clients having lower costs in the same time. Problem we would like to solve is lack of the tool that has simultaneously the following features:

- Free and Open Source
- Easy cloud management without complex knowledge
- Intuitive interface for individual workers to request resources
- Notifications for cost surpassing and redundant resources
- Manage notification well – only significant alerts
- Highly configurable and flexible
- Integration with communication via Email and Slack

2.2. Overview of existing solutions

Businesses that rely upon cloud services often reach a point where resources they are using up gradually become less and less manageable. As the problem is well known to the cloud market, both Amazon and other third-party companies made attempts to fulfill these needs by creating custom software fitting certain roles in optimizing AWS expenses that include:

1. Configuring budget limits and alerting users when they are exceeded
2. Instance alerting management
3. Cost analytics

Some of the most prominent tools currently available on the market that improve resource management experience for AWS cloud are described in the following sections. Every description is supposed to show key features crucial for cost optimization.

2.2.1. AWS Budgets

AWS Budgets[**AWSDocs**] is a part of Amazon Web Services that allows to set limits of a certain types that apply to a chosen period of time. When a limit is either exceeded, close to be exceeded or is forecasted to exceed the configured threshold before the end of that period, the administrator of that account is notified by email. Types of resources one can put this kind of a budget on include:

1. Money spent in total or on a certain type of machines.
2. Utilisation of selected services.
3. Utilisation or coverage of reserved instances.

2.2.2. AWS Cost Explorer and AWS Cost Management

AWS Cost Explorer[**CostExplorer**] enables access to all budget data. User can define and generate custom reports in a form of a data chart spanning a selected time interval with chosen time granularity of the samples. AWS Cost Explorer is also a basis for AWS Cost Management, which is basically a set of predefined reports that form an easily accessible dashboard.

2.2.3. Cloudability

Atlassian's Cloudability[**Cloudability**] delivers a budget system functionality analogous to AWS Budgets along with tools for predicting future spendings and presenting the real cost of AWS resources in utilisation. In comparison to Amazon's native tools Cloudability also allows management of multiple accounts in the same time. It saves effort of having to set up budgets separately in every owned account.

2.2.4. Apptio

Apptio[**Apptio**] provides a set of tools that mainly focuses on analysis of expenses and their forecasts, managing them collaboratively, and planning future ones. They expose features that make it easier to discover underutilized resource, compare spendings with a database of similar benchmarks, organize resources into groups to make reports even clearer, and offer other useful management utilities.

2.2.5. Stax.io

Stax.io's [Stax.io] main focus is to provide insight about cost, wastage, compliance and cloud quality. It can analyse how cloud resources are used, measure quality of the way cloud is utilized, set up checks for business-compliance of our cloud with several standards and give customized advice on what could be optimized to reduce wastage, while also allowing for creation of custom views of data. Basic tools for budgeting instances, accounts, tags and more, monthly or annually, and configuring overspend alerts are also available there.

2.2.6. SnowSoftware

SnowSoftware's toolset, alongside fulfilling some of the more specific usecases like optimizing usage of software from SAP Software Solutions or optimizing and managing software licenses, also has tools that are dedicated to optimizing cloud costs.

Snow for SaaS attempts to give a holistic view about application usage including, among others, how SaaS applications are used on cloud and whether there are zombie virtual machines [SnowSaaS].

Snow Automation Platform suggests approach based on automated and preconfigured provision of resources. By pre-giving those resources a decommission date one can avoid issue of zombie instances. It is also possible to preconfigure budgets and schedule machine starts and stops to further optimize costs [SnowBlog].

2.2.7. Conclusions

Using resources available to us, we concluded that most of the competing tools available on the market fail to provide both instance budgeting and machine termination automation. A number of them also puts most of their emphasis on analysing usage data rather than helping with instance management. From tools listed above, Snow Automation Platform is the only one that features both budgeting and automating machine termination and stopping. However, it's approach is not to straightforwardly manage already existing machines, but rather to automate their provisioning.

2.3. Our Solution

TODO

Chapter 3

Project development

In this section we will describe the way we developed the application – managment and organization of work along with supporting tools.

TODO

3.1. Version control

During development as our version control system we use Git and whole code repository is hosted on GitHub. We have chosen this particular service because of a few reasons.

First of all GitHub[**GitHub**] is one of the most popular Git hosting service and also has one of the biggest open source communities. It is quite important if we consider possibility of future development of our tool even after finishing this thesis.

The second reason is integration with lots of useful applications. It is easy to use it with different services improving development process like TravisCI[**TravisCI**] and Slack[**Slack**].

Last but not least, GitHub has built in issue tracker. It can be highly customised and makes project management easier. We created dashboard divided into a few sections, which separate tasks in different stages of development. Everything is automated – depending on actions performed by users like Pull Requests issues are moved between sections. In order to help prioritize tasks we introduced issues tagging system.

3.2. Continuous Integration

For Continuous Integration we use hosted service – TravisCI. It allows us to make sure new changes we introduce are compliant with the rest of our codebase. There are two reasons we have chosen this service: it is easy to configure using only one YAML file inside repository and it is free for open source projects. After every code submission, at the beginning TravisCI performs Smoke Tests trying to build the project and check whether it runs or not. Successful build is followed by both backend and frontend tests.

3.3. Communication

All of the communication happens via Slack. It is a platform allowing team members to communicate without use of email or group SMS. Various channels can be created there, with persistent messages history both public and private. There is also possibility to reach out to every member directly if needed. We use Slack due to the fact it is free and integrates easily with GitHub.

3.4. Organization of work

We work in iterations that last around one month. At the beginning we define and create new issues. During the iteration everyone chooses desirable task, completes it and makes corresponding Pull Request in GitHub. At the end we discuss our overall progress and plan next steps.

In order to let a new feature become part of the repository it has to be accepted. It means that Pull Request has to pass Continuous Integration system build and be approved by one of the reviewers. Every issue is solved in a separate branch and every PR is merged directly to the master branch, which allows us to group all commits that are part of one feature or issue.

3.5. Contact with client company – Sumo Logic

TODO (not sure if right title)

Chapter 4

Tool for AWS cost optimization

In this chapter we will describe the whole application – *Omigost*.

4.1. Overall architecture

Our solution is based on an architecture that is visualized on Figure 4.1

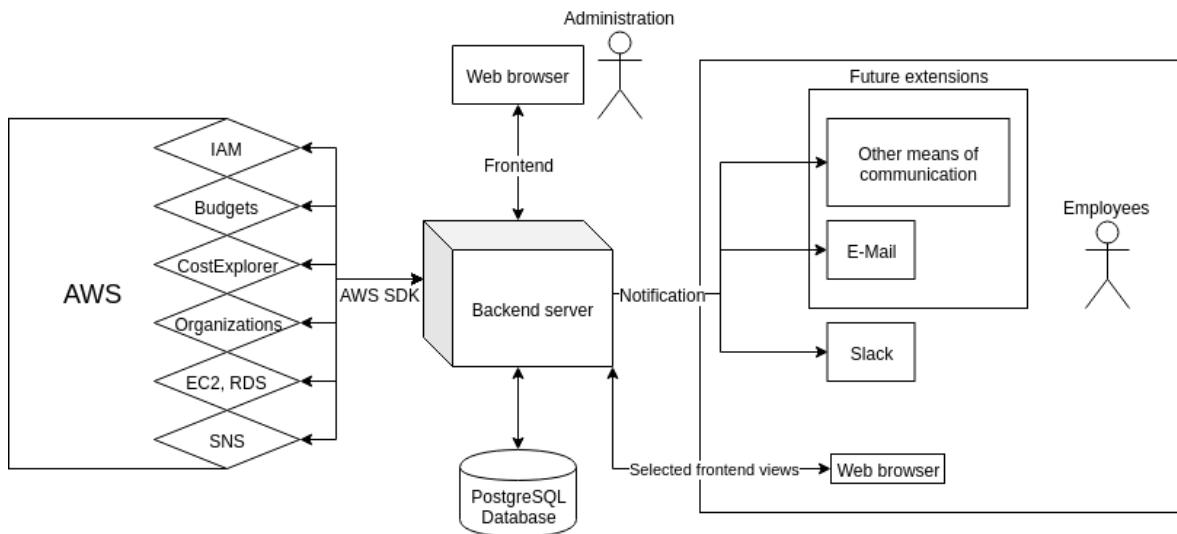


Figure 4.1: Draft diagram of the tool architecture

The architecture of our tool can be described as a traditional client-server architecture. Actors make requests over the HTTP protocol to the system with a web browser using the frontend client application served by the Omigost server or, alternatively, from any kind of a storage service like, for example, Amazon’s S3. Those requests are then handled by a backend application running on a central server machine. During the handling process our application uses both a local instance of a PostgreSQL relational database and a connection with selected Amazon Web Services via a SDK library. Additionally, end users are notified about budget or instance events related to their cloud activity through Slack.

4.1.1. Role of AWS cloud services

Connecting with Amazon's cloud services not only allows the application to access instance and spending data that is crucial to be able to fulfill the business requirements, but many of them also solve many of the problems we would otherwise have to solve ourselves – resulting in time saved and a codebase that is more concise and easier to maintain.

Services that are used in the project and their respective roles in it are as follows:

- Identity and Access Management (IAM) – provision of the main Amazon account ID that is used in other services,
- Organizations – insight into structure of the Organization used by the company, including listing accounts,
- Budgets – configurable spending watches that trigger a SNS notification whenever a budget limit is or is going to be exceeded,
- SNS – alert notifications for backend that let it know about spendings events,
- Cost Explorer – data for graph visualizing spendings in frontend,
- EC2/RDS – information about state of instances running on the account.

TODO at the end of the project – check if any other "sublibraries" were used and/or if any other use cases emerged

4.1.2. Backend

Backend is the main agent of our application's functionalities. It is responsible for fetching or receiving data from either of other connected resources, parsing it and taking appropriate actions. Some of its more important tasks include:

- receiving SNS budget alert notifications and notifying appropriate users via Slack,
 - providing single use tokens and links for non-admin users to be able to request budget limit increases,
- fetching and providing raw data for the frontend so it can be displayed to the user,
- receiving configuration or data modification requests and adjusting the application environment components to fulfill them,
- checking machine state on preconfigured times and notifying users about possible instance optimizations.

4.1.3. Frontend

Frontend is a web application that allows the user to see and alter the state of the application using a graphical user interface (GUI). Its main roles are to:

- request and parse raw data from backend into visual representations of it,
- provide an interface for the user,
- translate the interface clicks to appropriate backend requests and display the results of these requests.

4.2. Technology stack

4.2.1. Frontend

TODO

4.2.2. Backend

Core of the backend is *Java 8*. We have selected this programming language because most of us were familiar with it. Also this language has huge community and lots of useful libraries and frameworks.

For creation of the whole backend service we have chosen Spring [**Spring**] framework. It is currently one of the most common choices. Spring allows to build web applications imposing usage of design patterns like Model-View-Controller and Dependency Injection. It helps to keep code easily testable and well organized. Additionally it provides various features out of the box, which speeds up development process.

Gradle [**Gradle**] is responsible for build and dependencies management. It is easy to use with various IDEs and has flexible configuration. The main competitor of Gradle is Maven so obviously we considered both. Gradle turned out to be better in terms of performance and usage convenience.

For data persistence we have chosen PostgreSQL database. Relational database enforces having strictly defined data model with validation. Performance is not an important aspect for us so there was no need for other types of databases. PostgreSQL is a great open source database with strong data integrity and fault-tolerance guaranties. Furthermore, AWS lets set up PostgreSQL instance with just few clicks with automatically configured parameters for optimal performance.

In our application we use the following libraries:

1. Spring Core, Spring Boot, Spring Data, Spring Web
2. Project Lombok
3. JUnit
4. AWS SDKs

4.2.3. Deployment

The whole tool is bundled up by Docker. Following a straightforward instruction everyone is able to create their own Docker image of the application. Moreover, all of the configuration is present in one file making it easy to adapt it your way. A preferable method of deployment is to use AWS Elastic Beanstalk with Docker image. In such a way AWS will be responsible for almost everything including capacity provisioning, load balancing and auto-scaling.

4.3. Data models

TODO

4.4. Views and design

TODO

4.5. API

TODO

4.6. Communication integration

TODO

4.7. Configuration

TODO

Bibliography

- [Apptio] <https://www.apptio.com/products>.
- [Armbrust] 2009 M. Armbrust, A. Fox, R. Griffith, A. D. Joseph, R. H. Katz et al. *Above the Clouds: A Berkeley View of Cloud Computing*. Technical Report No. UCB/EECS-2009-28. <http://www.eecs.berkeley.edu/Pubs/TechRpts/2009/EECS-2009-28.html>.
- [AWSCostManagement] <https://docs.aws.amazon.com/aws-cost-management/latest/APIReference/Welcome.html/>.
- [AWSDocs] Amazon, *Amazon AWS documentation*.
- [Cloudability] <https://www.cloudability.com/product/plan/>.
- [CloudabilityAlerts] <https://blog.cloudability.com/creating-budget-alerts-by-tag-with-cloudability/>.
- [CostExplorer] <https://aws.amazon.com/aws-cost-management/aws-cost-explorer/>.
- [Docker] <https://www.docker.com/>.
- [GitHub] <https://github.com/>.
- [Laatikainen] 2013 G. Laatikainen, A. Ojala, O. Mazhelis *Cloud Services Pricing Models*.
- [Markus] 2011 M. Böhm, S. Leimeisier, C. Riedl, H. Krcmar *Cloud Computing and Computing Evolution* Technische Universität München (TUM), Germany.
- [Project Lombok] <https://projectlombok.org/>.
- [Slack] <https://slack.com/>.
- [SnowBlog] <https://www.snowsoftware.com/es/blog/2018/06/16/true-cost-aws>.
- [SnowSaaS] <https://www.snowsoftware.com/int/snow-saas>.
- [Spring] <https://spring.io/>.
- [Stax.io] <https://www.getapp.com/it-management-software/a/stax/>.
<https://www.stax.io/features>.
- [TravisCI] <https://travis-ci.org/>.