

Feedback Driven Development of Cloud Applications

Feedback Driven Development of Cloud Applications

Master-Thesis von Harini Gunabalan aus

Tag der Einreichung:

1. Gutachten: Prof. Dr.-Ing. Mira Mezini
2. Gutachten: Dr. Guido Salvaneschi, Dr. Gerald Junkermann, Aryan Dadashi



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Department of Computer Science
Software Technology Group

Feedback Driven Development of Cloud Applications
Feedback Driven Development of Cloud Applications

Vorgelegte Master-Thesis von Harini Gunabalan aus

1. Gutachten: Prof. Dr.-Ing. Mira Mezini
2. Gutachten: Dr. Guido Salvaneschi, Dr. Gerald Junkermann, Aryan Dadashi

Tag der Einreichung:

Bitte zitieren Sie dieses Dokument als:

URN: urn:nbn:de:tuda-tuprints-12345

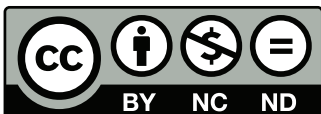
URL: <http://tuprints.ulb.tu-darmstadt.de/1234>

Dieses Dokument wird bereitgestellt von tuprints,

E-Publishing-Service der TU Darmstadt

<http://tuprints.ulb.tu-darmstadt.de>

tuprints@ulb.tu-darmstadt.de



Die Veröffentlichung steht unter folgender Creative Commons Lizenz:

Namensnennung – Keine kommerzielle Nutzung – Keine Bearbeitung 2.0 Deutschland

<http://creativecommons.org/licenses/by-nc-nd/2.0/de/>

Erklärung zur Master-Thesis

Hiermit versichere ich, die vorliegende Master-Thesis ohne Hilfe Dritter nur mit den angegebenen Quellen und Hilfsmitteln angefertigt zu haben. Alle Stellen, die aus Quellen entnommen wurden, sind als solche kenntlich gemacht. Diese Arbeit hat in gleicher oder ähnlicher Form noch keiner Prüfungsbehörde vorgelegen.

Darmstadt, den 10th May 2016

(Harini Gunabalan)

Abstract

Over the last few years, the Cloud Computing Paradigm has gained a lot of importance in both the Academia and the Industry. The cloud has not only changed the IT Landscape from the user's perspective but has also changed how the Developers develop applications on the cloud. The increasing adoption of the DevOps approach has led to the removal of the boundaries between the development and the operations.

Among the the three levels of the Cloud Computing: Infrastructure-as-a-Service (IaaS), Platform-as-a-Service(PaaS) and Software-as-a-Service(SaaS), Software Developers are mainly concerned with the PaaS which allows them to focus on the Application Development. Leveraging the fact that the Application is hosted on the cloud, there are additional metrics regarding the application available to Developers in different Log formats. There are also Cloud Monitoring Tools which consolidates these logs with a huge volume of data representing the run-time metrics etc. of the cloud applications. Though monitoring of the cloud applications has been done by many tools, most of the developers do not go through the cumbersome error/warning log data. These log data are not visually made available to the developers in their Development Environment. By providing this information, developers can have an overview of the Application Performance at real time, and can capture issues that occur at scale which normally could not be captured by Profilers. This could also help us to reduce the developer-operator gap providing an improved DevOps experience.

This Thesis work aims in addressing this issue. The focus is mainly to map the run-time metrics to the source code artifacts, thereby helping the developers link the Run time metrics to the source code. Mapping will involve log data aggregation, code analysis techniques etc.

Contents

A	Introduction	5
A.1	Motivation	5
A.2	Problem Statement	6
A.2.1	Continuous Delivery in Cloud Development	7
A.2.2	Monitoring Metrics: Scalability and Availability of the Application	7
A.2.3	Incorporating the collected Feedback in the Development Environment	7
A.3	Contribution	8
A.4	Structure of the Thesis	8
B	State of the Art	9
B.1	Cloud Monitoring	9
B.2	Auto-scaling	10
B.3	Data Modeling	12
B.3.1	State Space Models	13
B.3.2	Polynomial Models	14
B.4	Performance Analysis using Source Code History in Evolving Software	15
B.5	Feedback Driven Development	17
C	System Design	19
C.1	System Architecture Overview	19
C.2	Monitoring Metrics	19
C.3	System Architecture Details	19
C.4	Phase 1: Modeling	20
C.5	Phase 2: Evaluation of the Model	20
D	System Implementation	21
D.1	System Implementation	21
E	Evaluation and Results	22
E.1	Evaluation	22
F	Conclusion and Future Work	23
F.1	Conclusion and Future Work	23

List of Figures

1	Cloud Computing Stack	6
2	Cloud Monitoring [1]	10
3	Architecture of polyglot autoscaler [2]	12
4	Positive, Negative and no Correlation between X and Y	13
5	System represented by State Space Model	14
6	Source Code Changes of two versions in Agilefant [4]	16
7	Code Changes that caused maximum Performance Variations [5]	17

A Introduction

Cloud Computing is one of the fields in Computer Science that has gained rapid growth and importance in the recent years. The fact that the servers are remote hosted rather than local servers has led to innumerable small scale businesses. Start-ups no longer require high infrastructure, instead they just need to pay for the amount of resources that they actually use (pay-per-use). This also significantly reduces the initial monetary setup costs for such start-ups. There has been extensive research in Cloud Computing areas such as auto-scaling of resources at the infrastructure level, monitoring of metrics at both the application and infrastructure level. However, there is not much research done in how these monitored metrics are utilized by the Cloud Application Developers. Making the run-time metrics visibly effective for the developers in their Development Environment is the issue this thesis is aiming to solve.

This chapter is structured into four sections. The first section provides the motivation of this thesis work. The second section gives an insight into the problem statement which this Thesis work aims to solve. The third section details the contribution of the work and the fourth section provides how the following chapters of the Thesis are structured.

A.1 Motivation

Software Engineering Practice in the industry has faced a phenomenal change since the advent of Cloud Computing. This is mainly due to the flexibility and the dynamic scaling up and scaling down of the infrastructure as required by the current workload. This proves not only to be elastic but cost-effective as well. It is quite obvious that this elasticity is achieved by the continuous monitoring of several metrics that indicate the demand at the moment, and provisioning the necessary resources to meet the monitored demand. Distributed, scalable enterprise-wide applications also mandate the monitoring of metrics for reasoning the effectiveness of the applications by engineers and business analysts [6].

The metrics that are being monitored vary widely. For instance, the metrics such as memory consumption, CPU utilization, Network bandwidth utilization could be considered to be at the Infrastructure level, whereas some other metrics such as Response times of methods/procedures, the number of users accessing the application, maximum number of users who can use the application simultaneously, etc. could be considered to be at the application level. Sometimes, the application level metrics could depend on the primitive metrics at the infrastructure level and vice-versa.

As Cloud Application Development has becoming more common, the run time monitoring metrics of these applications are available through several Application Performance Monitoring (APM) tools such as Amazon Cloudwatch, New Relic, etc. But, they do not provide any valuable and visible feedback to the developers, and hence most of the cloud developers do not use it. However, this run-time monitoring data could be used to provide useful analytic information such as performance hotspots that is taking a lot of execution time, and predictive information such as methods or loops that may become critical, even before the deployment. This type of analytic and predictive feedback should be provided to the developers in their IDEs which otherwise may not be explored by the developers. This technique of utilizing the monitoring

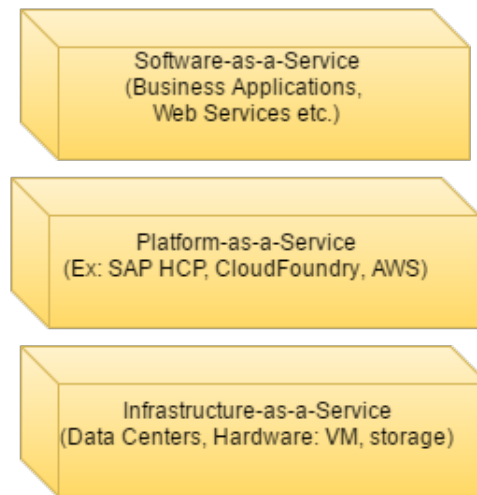


Figure 1: Cloud Computing Stack

data is known Feedback driven development. Feedback Driven Development that provides visual tools to the cloud developers is the focus of this thesis research.

A.2 Problem Statement

The Cloud computing paradigms are classified into three service models which forms a stack as shown in Figure 1.

[7] defines a framework that collects data as defined by a 3-D cloud monitoring model. From the software engineering perspective, it is important to note how Cloud Computing impacts the development practices. Based on the research conducted, there are two important research issues [8].

- Impact of DevOps on Cloud Application Developers: The Cloud Developers are forced to look into the huge log data of the cloud applications. Development and Operations are brought much closer than the traditional software methods. Sometimes the same person acts as both the developer and the operator.
- Data and Tools utilized by Cloud Developers: The data produced by the logs include Business Metrics, System-level Data etc. Sometimes the implementation changes may even have monetary consequences in the cloud, nevertheless most of the developers do not pay attention to this aspect. Hence it can be argued that once these operational metrics are brought closer to the Cloud Developers Environment, they would be able to pay closer attention to this information thereby achieving cost-effective applications.

Considering these issues, it is certainly important how we are going to leverage these Cloud monitoring Logs to make a useful impact for the Cloud Developers.

A.2.1 Continuous Delivery in Cloud Development

As we compare the software development on-premise with that of the Cloud, there has been a huge change in every software release. Deployment cycles has been reduced from months to days and sometimes even within a few hours the next version is released. This process is often referred to Continuous Delivery(CD) in the Cloud Computing terms. Most companies make use of this Continuous Delivery to rollout new features and evaluate their new ideas in a controlled manner [9]. CD has become a huge success and companies such as Google,Facebook etc. adopt CD of varying degrees for some of their services. When a feature is delayed for the current roll-out, it gets delayed by months when the traditional software development life-cycle is used. The feature needs to be delayed until the next release. Whereas in this CD approach the next release could be in the same day or in the same week, leading to small changes of production code. This also leads to a state called "perpetual development" where the code is always under continuous development and there is no stable release version for a particular product.

Owing to this new release paradigm, there are a lot of extra information generated. The live performance of the application, click-streams from the User Interface of application, error and warning logs, infrastructure related data etc. are produced. There are a lot of existing APM Tools that collect this data and generate information out of it, however, how this information could be made effective to the software developers in their daily routine is a topic that is not discussed very often.

A.2.2 Monitoring Metrics: Scalability and Availability of the Application

Some of the monitoring metrics collected include the CPU usage, response time of the request, number of instances the application is hosted on, no. of requests each instance serves during a particular time period, error logs of the request etc. While these metrics focus mainly at the Infrastructure level, the logs instrumented into the application are also collected. We use these metrics to focus on two main challenges of cloud computing: Scalability and Availability of the Cloud applications.

Scalability is the ability to increase or decrease the resources of an instance or the number of instances so that the changing demand of the incoming requests can be met. The platform related log data are collected for both Horizontal scaling (increasing the number of instances) and Vertical scaling(increasing the resources of each instance).

Availability is one of the major goals of Cloud Computing. It means that the Cloud Services need to be available and accessible at anytime from anywhere. For business to happen continuously, it is necessary for the services to be highly available. The definition of availability is specified by the different Cloud vendors in their SLAs. For example, Google Search is known for its high availability.

A.2.3 Incorporating the collected Feedback in the Development Environment

Using the collected metrics, in this thesis, we aim to provide an efficient feedback to the Cloud Developers. The collected feedback is integrated into the Development Environment (IDEs) so

that developers are able to utilize the feedback to make their applications better scalable and highly available.

A.3 Contribution

A.4 Structure of the Thesis

This thesis is structured in six chapters. Chapter 1 provides a short introduction into the topic and describes the goals of the thesis. Chapter 2 includes more background information and presents the state of the art of the important topics of this thesis: Cloud Application Performance Monitoring(AMP) tools, Feedback Driven Development(FDD) in general, and how FDD could be useful for Cloud Application Developers. The third Chapter describes an overview of the high-level System Design and the design decisions made in this research work. The fourth chapter explains the system on a lower fine-grained level. Interesting Implementation details are also provided here. Chapter 5 shows the various case studies, evaluates the developed system and illustrates its usage as well as possible applications. Finally, Chapter six provides the conclusion of the thesis and outlines the future work ideas.

B State of the Art

This chapter presents the state of art of the topics relevant for this thesis. In the first section, we explain what is Cloud Monitoring. The second section briefs about Auto-scalers and third section provides a background on Data Modeling. The fourth section describes about mining source code changes to identify performance regressions and the final section details about Feedback Driven Development and the types of FDD.

B.1 Cloud Monitoring

As Cloud Computing is gaining popularity, the need for Cloud monitoring is becoming increasingly important to both the Cloud Providers and the Cloud Consumers. At the Cloud Provider side, Cloud Monitoring is the key principle behind which the actual controlling of hardware takes place. It enables them to scale the infrastructure, if necessary. On the Cloud Consumer side, Cloud Monitoring enables to check the Availability, QoS etc. of the applications. The consumers can verify any SLA violations by comparing the Key Performance Indicators(KPI) parameters provided by Cloud Monitoring.

[1] explains in detail about the need for Cloud Monitoring: the basic concepts involved, the properties which needs to be maintained, and finally also lists down the open issues with respect to Cloud monitoring. These are summarized in Figure 2.

There are several Cloud monitoring platforms and services such as CloudWatch [10] [11], AzureWatch [12] , NewRelic [13] etc. Table 1 provides a list of Cloud Monitoring platforms and services. Amazon CloudWatch provides users the monitored information for 2 weeks. Users are allowed to plot these information, set thresholds, alerts etc and these alerts can be used to perform any substantial action such as sending an Email or even in AutoScaling [14]. AutoScaling is explored in detail in the next section.

Table 1: Cloud Monitoring Platforms and Services

Cloud Monitoring Platforms	Cloud Monitoring Services
CloudWatch [10] [11]	New Relic [13]
Nimsoft [15]	Cloudyn [16]
AzureWatch [12]	Up.time [17]
Nagios [18]	CloudSleuth [19]
Nimbus [20]	Cloudstone [21]
GroundWork [22]	Boundary [23]
LogicMonitor [24]	Cloudfloor [25]
CloudKick [26]	CloudClimate [27]
Monitis [28]	CloudHarmony [29]

Some other Cloud Monitoring Platforms such as Nimsoft Monitoring Solution [15] provides a unified monitoring dashboard to view infrastructures provided by Salesforce, Rackspace, Google or Amazon. Nagios [18] is a popular open source Cloud Monitoring platform which provides monitoring of virtual machines and storage (Amazon EC2 and S3). It also supports OpenStack

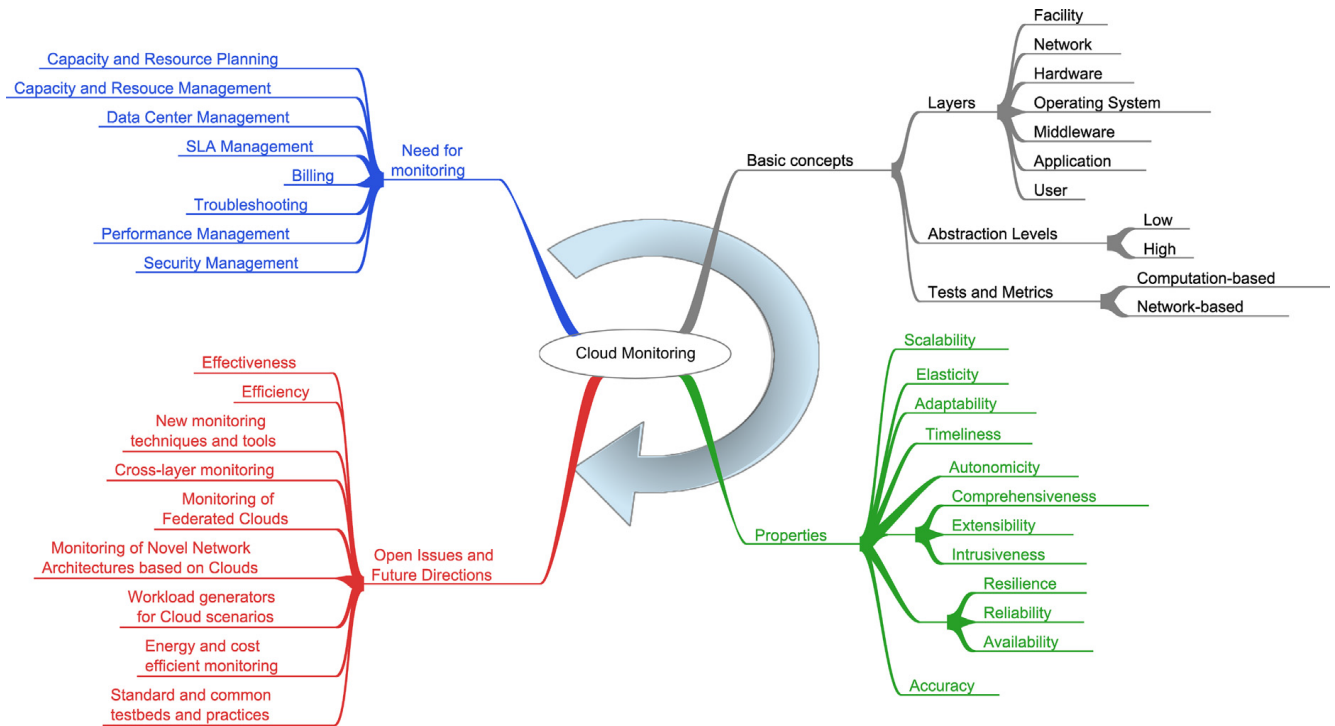


Figure 2: Cloud Monitoring [1]

[30], an open Source Cloud IaaS. New Relic [13] is a web-based Monitoring service that helps to monitor the application infrastructure and performance, adhering to timeliness, resilience, availability and accuracy.

While all of the above Cloud Monitoring platforms and services are great, most of them do not consider multiple layers or real-time data. Some of them consider multiple layers whereas do not take into account the real-time data, and some consider real-time data but do not consider the multiple cloud layers [31]. [7] proposes a 3D-Cloud Monitoring framework called the Ceiloesper framework, which combines monitoring in multiple layers with real time data and it also performs the data analysis for multiple management actions. It is based on the Complex event processing (CEP) and uses the Esper CEP Engine.

B.2 Auto-scaling

Scaling of Cloud Infrastructure means changing the current infrastructure. It could be of two types: Horizontal and Vertical Scaling. Horizontal Scaling is a methodology of adding/removing machines whereas vertical scaling is increasing/decreasing the resources such as CPU/Memory/Disk to existing machines.

The Metrics monitoring and Data Collection has a huge importance with respect to Auto-scaling. Auto-scaling or automatic scaling is a process where the Cloud Platform adapts itself by increasing or shutting down the number of instances on which the application is currently deployed depending on the current load. For enterprises running their own infrastructure shutting down servers that are not being utilized could save electricity costs, whereas for enterprises

running their applications on Cloud, auto-scaling could lead to saving costs due to the pay-per-use model of the Cloud. Auto-scaling also improves the efficiency of applications. In the scenario mentioned in [32], auto-scaling improves the instance utilization of the open source AppScale PaaS by 91% and it also brings down the average time taken to serve the requests. Auto-scalers can be broadly classified as the following:

- **Reactive Auto-scaling:** Auto-scaling as provided by most of the cloud providers such as Amazon Web Services, Microsoft Azure, and IBM Bluemix etc are reactive. This means based on monitoring the relevant metrics, whenever a certain metric increases or decreases beyond a particular predefined threshold, additional instances are added or removed. This method which is more of a rule-based mechanism is a reactive auto-scaling method. This is easier to be implemented as it involves monitoring metrics, and framing rules and policies for scaling. While this method serves in most scenarios, the question arises whether it is capable to handle bursty traffic.

In [2] the rule based reactive autoscaler of IBM's Bluemix PaaS, Polyglot application, is described. Polyglot autoscaler allows application developers to set thresholds based on which instances need to be added (scale-out) or removed (scale-in). These threshold values could be parameters such as CPU Utilization, memory and heap usage. The architecture of the Polyglot autoscaler can be found in Figure 3. It consists of the four components: Agents which collect the performance information, a monitoring service which continuously monitors the health of the cloud application, a scaling service which makes the decision of whether scaling needs to be performed or not, and a Persistence service to keep track of the enactment points (points where the application is scaled in time).

- **Predictive Auto-scaling:** Predictive auto-scaling comes very handy to handle bursty workloads. By analyzing the historic time series data, it may be possible to predict the workload at a future time, thereby enabling predictive auto-scaling. The effectiveness of this method depends on the efficiency of the workload prediction.

[33] introduces a predictive auto scaling technique that uses a Machine Learning engine to make predictions based on a deadline driven algorithm for predicting the future state of the system. [34] and [35] describes a predictive auto-scaling tool, Scyer, used by Netflix to provision the correct number of Amazon Web Services [36] instances. This is different from the Amazon AutoScaling(AAS) [14], which is a reactive one. Scyer's prediction engine is able to provision the resources based on two prediction algorithms to predict the workload. The prediction algorithms implemented are augmented linear regression based algorithm and Fast Fourier Transformation based algorithm.

- **Hybrid Auto-scaling:** Hybrid Auto-scaling is a combination of both the Reactive and Predictive approaches. As mentioned in [34], Scyer tool works in co-ordination with the AAS for more efficient auto scaling. [37] describes the architecture and implementation of Platform Insights, which is another hybrid auto-scaler that employs a reactive rule-based and a predictive model-based approach in a coordinated manner.

Auto-scalers face the following problems as listed in [38]:

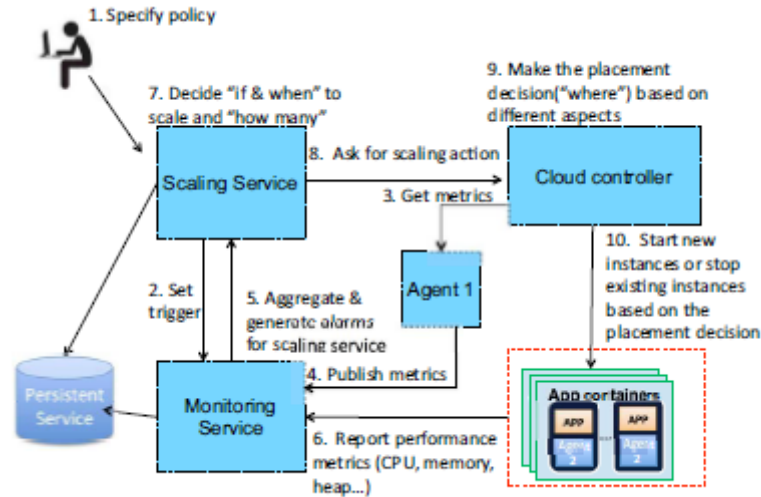


Figure 3: Architecture of polyglot autoscaler [2]

Under Provisioning: The application is hosted on lesser infrastructure than that is necessary to process all the incoming requests. Due to Service-level-agreements (SLA's), it takes a while for it to reach up to the required amount of infrastructure. This could also lead to SLA violations.

Over Provisioning: There is no SLA violations in this scenario. However the actual amount of resources is greater than the required amount of resources and hence the customer could be paying extra cost that his actual usage.

Oscillation: When there is an oscillation between Under Provisioning and Over Provisioning it causes an undesirable and unstable state.

In order to solve this, auto scaling needs to focus on the MAPE Loop [38]: Monitor, Analyze, Plan and Execute. The necessary monitoring metrics are collected and analyzed to decide on the type of autoscaling: reactive/predictive/hybrid. The planning phase is done on how to actually perform the scaling: Horizontal/Vertical. Finally the actual scaling is performed based on SLA.

B.3 Data Modeling

Correlation and Covariance are two important concepts in statistics. Both indicate how closely two variables are related. For instance if variable X increases, will variable Y increase or decrease or does not depend on X. In addition to this, correlation helps us to understand to what extent the two variables change with respect to each other.

Correlation could be either positive or negative. If the variable Y increases proportionately when variable X is increased by a unit, it is known as positive correlation. On the other hand if the variable Y decreases proportionately when variable X is increased it is a negative correlation. This can be explained graphically as shown in Figure 4. If all the points are centered around the straight line: $Y = X$, then X and Y are said to be positively correlated. Whereas if all the points

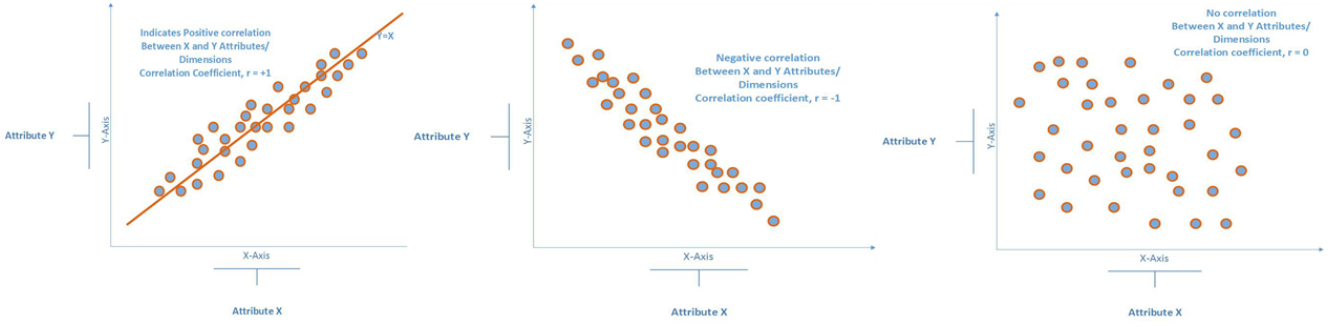


Figure 4: Positive, Negative and no Correlation between X and Y

are centered around a line $Y = -X$, then X and Y are said to be inversely correlated. If all the points are scattered throughout then there is no correlation between variables X and Y . Figure 4 depicts positive, negative and no correlation.

According to statistics, if X_i and Y_i are sample data for the two variable under consideration then correlation can be calculated as [3]:

$$\text{Correlation, } r_{xy} = S_{xy} / S_x S_y$$

where S_x = sample standard deviation of variable X , S_y = sample standard deviation of variable Y and S_{xy} is the sample covariance of the variables X and Y . The correlation coefficient values r_{xy} ranges between -1 and 1. If the value is positive, then there is a positive correlation and if the value is negative, then there is a negative correlation. If the value is 0 then there is no correlation. Also, the closer the values are to +1 or -1, the stronger is the correlation, positive or negative respectively.

So far we considered only a single input variable, X and a single output variable, Y . However in reality most of the systems tend to be multiple-input multiple-output (MIMO) systems rather than the single-input single-output (SISO) system. The data that we deal in real world does not contain just 2 attributes. Most of the real world scenario involves a minimum of 5 to 6 dimensions and depending on the applications this may go as high as 20 or even more. Hence we explore further into multivariate correlation models: State Space Models, and Polynomial Models.

B.3.1 State Space Models

State space model represents a system by a set of First order differential equations, and state variables. Mathematically it can be described that the output $Y(t)$ of a system at time, t can be predicted for any time $t > t_0$, where t_0 is an initial time, and provided that we know the input and output of the system at time t_0 and a minimum set of variable $x_i(t)$ where $i = 1$ to n . In this case, n is the order of the state space model [39].

Figure 5 shows a system described by a state space model. The vector $u_1, u_2, u_3, \dots, u_i$ are the inputs while the output vector is y_1, y_2, \dots, y_k . By knowing the inputs and outputs at time t_0 the state variables: $x_1, x_2, x_3, \dots, x_n$ are first measured. Then it becomes possible to predict the output at any future time, t by knowing the inputs at that time and the measured state variables.

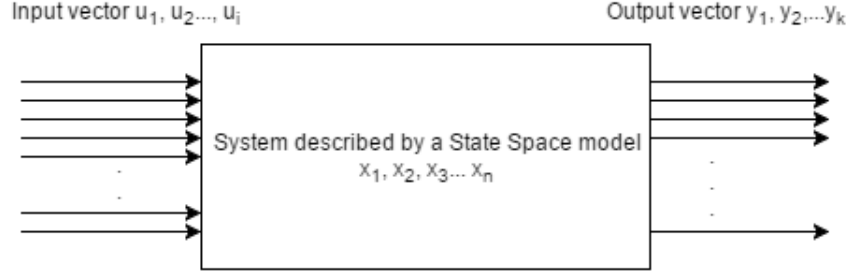


Figure 5: System represented by State Space Model

In state space modeling, the time derivative of the state variables are represented as a function of the state variable, and inputs, $dx/dt = f(x, u, t)$. Considering a Linear Time Invariant (LTI) systems, we have the state equation [39]:

$$dx/dt = Ax + Bu$$

where A and B are matrices with constant coefficients that weight the system's state variable and inputs respectively. Similarly the output equation can be written as [39]:

$$y = Cx + Du$$

where C and D are matrices with constant coefficients that weight the system's state variables and inputs respectively. There are several physical systems where the D matrix is found to be a null matrix thereby reducing the output equation to $y = Cx$, where the output depends on a weighted combination of the state variables.

B.3.2 Polynomial Models

Some physical systems do not always adhere to linear equations. Hence to model these type of systems we consider polynomials. Additionally there could be systems which depends on the previous values of the inputs, previous values of the outputs as well. Based on these we have the following four polynomial models [40]:

1. ARX Model: The ARX model to evaluate the output is based on Auto-regression (the past output values) and inputs. Auto regressive model is a model whose current output depends on the past values. The generic notion to denote auto-regressive model of order p, AR(p) for a variable X is:

$$X_t = c + \sum_{i=1}^p \rho_i X_{t-i} + e(t)$$

where c and ρ_i are constants and e(t) is the noise [41]. Considering auto regression and the inputs the ARX model can be mathematically described as:

$$A(z)y(t) = B(z)u(t - n) + e(t)$$

where y(t) is the output, u(t) is the input, and e(t) is the noise/error measured in the output. A(z) and B(z) are polynomials of the specified order with respect to the backward shift operator Z^{-1} . For example, $Z^{-n}u(k) = u(k - n)$ [42].

2. ARMAX Model: Unlike the ARX model, in ARMAX the stochastic dynamics are considered. Therefore this model handles a system where there is a domination of noise. ARMAX models are better for systems with more disturbances. In general, the moving average model of order q , MA(q) is represented in the below notation:

$$X_t = e(t) + \sum_{i=1}^q \theta_i e(t-i)$$

where θ_i are constants and $e(t)$ and $e(t-i)$ are the noise/errors [41]. The notation for the Auto-regressive moving average (ARMA) model is as below:

$$X_t = c + e(t) + \sum_{i=1}^p \rho_i X_{t-i} + \sum_{i=1}^q \theta_i e(t-i)$$

This model includes both AR(p) and MA(q) models. Based on these we have the following mathematical equation for the ARMAX model:

$$A(z)y(t) = B(z)u(t-n) + c(z)e(t)$$

where, $y(t)$ is the output, $u(t)$ is the input, and $e(t)$ is the noise. $A(z)$, $B(z)$ and $C(z)$ are polynomials of specified orders with respect to the backward shift operator Z^{-1} [43].

3. Output-Error Model: The notation for the Output Error model is as below:

$$y(t) = [B(z)/F(z)]u(t-n) + e(t)$$

where, $y(t)$ is the output, $u(t)$ is the input, and $e(t)$ is the noise. $B(z)$ and $F(z)$ are polynomials of specified orders with respect to the backward shift operator Z^{-1} [44].

4. Box-Jenkins Model: The notation for the Box Jenkins model is as below:

$$y(t) = [B(z)/F(z)]u(t-n) + [C(z)/D(z)]e(t)$$

where, $y(t)$ is the output, $u(t)$ is the input, and $e(t)$ is the noise. $B(z)$, $F(z)$, $C(z)$ and $D(z)$ are polynomials of specified orders with respect to the backward shift operator Z^{-1} [45].

B.4 Performance Analysis using Source Code History in Evolving Software

Software evolution is defined as the change of characteristics of a software with time. Continuous Delivery has led to continuously evolving software. This means frequent code changes are prevalent and this naturally causes performance regressions. Performance of a software is quite important and hence valuating performance regressions during code changes becomes a necessity. Performance regression can be defined as a state when the application under consideration behaves worse in a new code deployment compared to its previous deployment. In this section we look into two source code mining tools: PerfImpact [4] and LITO [5].

<p style="text-align: center;">V3.2</p> <pre> public String generateTree(){ Set<Integer> selectedBacklogIds = this.getSelectedBacklogs(); if(selectedBacklogIds == null selectedBacklogIds.size() == 0) { addActionError("No backlogs selected."); return Action.ERROR; } return Action.SUCCESS; } </pre>	<p style="text-align: center;">V3.5</p> <pre> public String generateTree(){ Set<Integer> selectedBacklogIds = this.getSelectedBacklogs(); if(selectedBacklogIds == null selectedBacklogIds.size() == 0) { Collection<Product> products = new ArrayList<Product>(); productBusiness.storeAllTimeSheets(products); for (Product product: products) { selectedBacklogIds.add(product.getId()); } } return Action.SUCCESS; } </pre>
<p style="text-align: center;">V3.2</p> <pre> public StoryTreeBranchMetrics calculateStoryTreeMetrics(Story story) { for(Story child : story.getChildren()) { StoryTreeBranchMetrics childMetrics = this.calculateStoryTreeMetrics(child); } return metrics; } </pre>	<p style="text-align: center;">V3.5</p> <pre> public StoryTreeBranchMetrics calculateStoryTreeMetrics(Story story) { for(Story child : story.getChildren()) { if (child.getId() == story.getId()) { continue; } StoryTreeBranchMetrics childMetrics = this.calculateStoryTreeMetrics(child); } return metrics; } </pre>

Figure 6: Source Code Changes of two versions in Agilefant [4]

1. PerfImpact: PerfImpact identifies the Performance Regressions and recommends potential code changes that has led to the performance degradation. PerfImpact achieves this as a two step-process:

- **Identification of Inputs which cause the Performance Regression:**
PerfImpact defines a "Fitness Function" that determines the inputs which cause the delay in execution of a newer code deployment V_{i+1} compared to its previous deployment V_i . The fitness function makes use of Genetic Algorithms to achieve this.
- **Mining execution traces to identify code changes that lead to Performance Regressions:**
PerfImpact also has a "Mining Function", which identifies those methods which take a longer execution time in V_{i+1} compared to V_i . These methods are tagged as Potentially problematic methods. Between the two deployments there could be several code changes/commits. Each code change is ranked based on the number of potentially problematic methods involved. The code changes with higher number of problematic methods are ranked higher and considered as the possible root cause for the performance regression.

PerfImpact was evaluated on two open source web applications: JPetStore [46] and Agilefant [47]. Figure 6 shows the source code changes in two versions of Agilefant. The evaluation shows that the inputs which cause Performance regressions are identified efficiently. PerfImpact also lists the potentially harmful code changes which could be used further in Code Inspectors and Root Cause Analysis.

2. LITO, a Horizontal Profiling technique: LITO is a cost model to determine if a code commit has caused performance regressions based on sampling the execution of versions. This approach resolves the following research questions (RQs) as below:

Source Code Changes		R	I	R/I	Total
1	Method call additions	23	0	1	24 (29%)
2	Method call swaps	15	9	0	24 (29%)
3	Method call deletion	0	14	0	14 (17%)
4	Complete method change	6	0	3	9 (11%)
5	Loop Addition	5	0	0	5 (6%)
6	Change object field value	2	0	0	2 (2%)
7	Conditional block addition	0	2	0	2 (2%)
8	Changing condition expression	0	2	0	2 (2%)
9	Change method call scope	1	0	0	1 (1%)
10	Changing method parameter	0	1	0	1 (1%)
Total		52	28	4	84 (100%)

Figure 7: Code Changes that caused maximum Performance Variations [5]

- RQ-1: Is there a set of specific methods which will cause performance variations when the source code of these methods are modified? According to [5], this is not really true. This is in contrast to PerfImpact. This approach was tested on 17 open source projects and the results showed that the methods, which cause performance variations before, not necessarily contributed to the performance variations in the newer versions.
- RQ-2 What are the recurring code changes which affects the performance of an evolving software? The major code changes the caused performance variations are method call addition, method call deletion, method call swap, Complete Method call change, and Loop Addition as compared to the other code changes listed in Figure 7 [5].

B.5 Feedback Driven Development

By analyzing any Cloud Application's logs, we can get hold of a huge amount of information. This can be broadly classified into: Application Level Logs and Infrastructure Level Logs. This data could be made useful to both the Developers and Operators. Infrastructure Logs provides details such as number of instances, the memory, CPU, Disk utilization of each instance, which instance serves a particular request etc. By making this kind of data visible to the developers, they can tweak the application development process, as they have access to the cloud internals. At the same time, the Cloud operators also benefit with the relevant business metrics to manage the instances more efficiently.

Collecting these run-time data, aggregating them into useful feedback, and feeding them back into the Development process of an application could create a useful impact in the future deployment of the application. This process is known as Feedback Driven Development(FDD). FDD can be classified into 2 types: Analytic FDD and Predictive FDD [48].

- Analytic Feedback Driven Development: Analytic FDD is the run time data from the previous deployments, which is brought directly into the developer environment. It provides a mapping between the log data collected and the source code artifacts. This helps the

developers to understand how the run-time metrics directly impact the source code. Based on this developers can alter and optimize the code based on the real time user behavior. In practice, Analytic FDD deals with visualizing the run time operations data and how it is being mapped to the code artifacts.

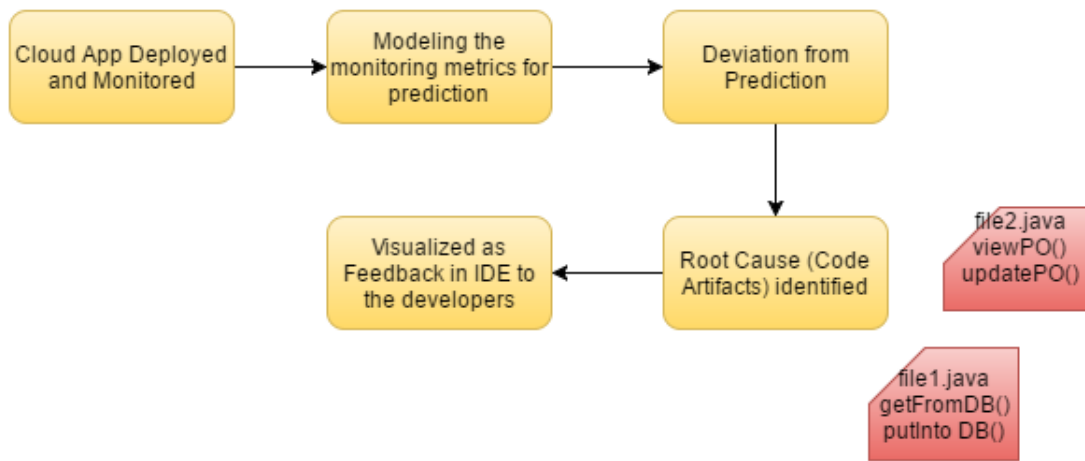
- Predictive Feedback Driven Development: Predictive FDD is one step ahead compared to the Analytic FDD. It utilizes the run-time feedback to warn the developers about the current code changes even before the updated source code is deployed. Predictive FDD is combined with static code analysis to give better predictions regarding a code change.

C System Design

In this chapter, the high level architectural design of the thesis work is explained. We look into the design decisions such as the choice of the metrics to be collected, the modeling methods and how the feedback is visualized in a productive manner to the developers.

C.1 System Architecture Overview

The architectural overview of the system is depicted in Figure C1. Initially, a sample application for demo purpose is considered. This application is deployed in the cloud systems and monitored continuously. The data collected during this monitoring is extracted into a clean format to be used for the Modeling.



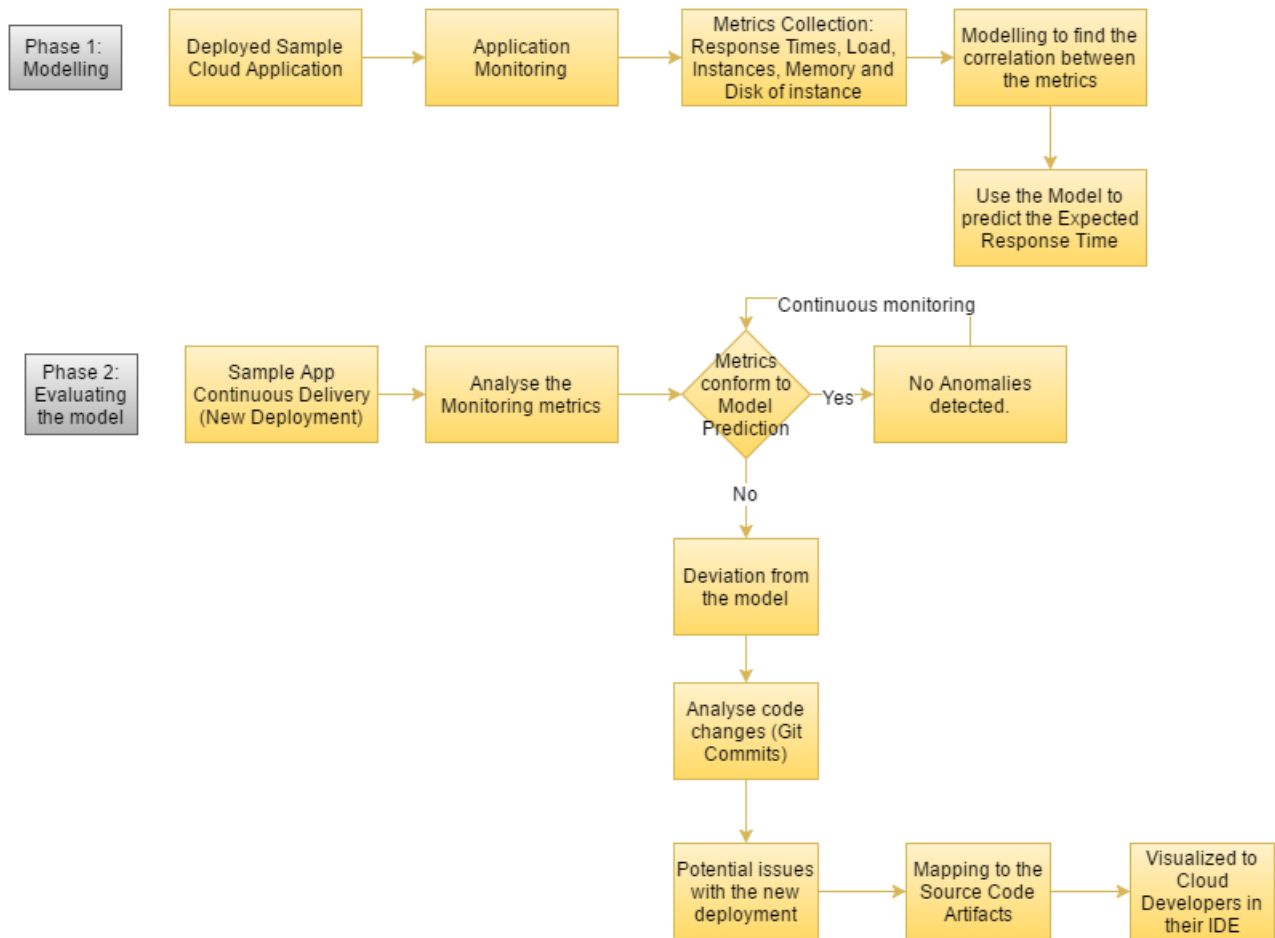
The MATLAB System Identification Toolbox is used to model the data obtained. Using the model, we estimate the expected performance of the sample application. If a deviation occurs from this estimation then the root cause of the issue is analyzed. This is done by investigating the code changes such as code commits. Finally, the source code artifacts that is related to the anomaly is identified and visualized in the developer's IDE. Owing to the increased use of Continuous Delivery and DevOps approach, providing immediate run time production scenario to the Development

C.2 Monitoring Metrics

The deployed cloud app is monitored continuously to collect the various dimensions. In this thesis work, we focus on the following metrics:

C.3 System Architecture Details

The overall system design is split into two phases as shown in the figure C2. The first phase explains the Modeling details and the second phase focuses on how the model is evaluated to identify the root cause that needs to be visualized as a feedback to the developers.



A sample application is deployed into the Cloud. This application should be available for monitoring the parameters.

C.4 Phase 1: Modeling

C.5 Phase 2: Evaluation of the Model

D System Implementation

D.1 System Implementation

E Evaluation and Results

E.1 Evaluation

F Conclusion and Future Work

F.1 Conclusion and Future Work

References

- [1] G. Aceto, A. Botta, W. De Donato, and A. Pescapè, “Cloud monitoring: A survey,” *Computer Networks*, vol. 57, no. 9, pp. 2093–2115, 2013.
- [2] S. R. Seelam, P. Dettori, P. Westerink, and B. B. Yang, “Polyglot application auto scaling service for platform as a service cloud,” in *Cloud Engineering (IC2E), 2015 IEEE International Conference on*, pp. 84–91, IEEE, 2015.
- [3] “correlation.” http://ci.columbia.edu/ci/premba_test/c0331/s7/s7_5.html.
- [4] Q. Luo, D. Poshyvanyk, and M. Grechanik, “Mining performance regression inducing code changes in evolving software,” 2016.
- [5] J. P. Sandoval Alcocer, A. Bergel, and M. T. Valente, “Learning from source code history to identify performance failures,” in *Proceedings of the 7th ACM/SPEC on International Conference on Performance Engineering*, pp. 37–48, ACM, 2016.
- [6] P. Leitner, C. Inzinger, W. Hummer, B. Satzger, and S. Dustdar, “Application-level performance monitoring of cloud services based on the complex event processing paradigm,” in *Service-Oriented Computing and Applications (SOCA), 2012 5th IEEE International Conference on*, pp. 1–8, IEEE, 2012.
- [7] D. Bruneo, F. Longo, and C. C. Marquezan, “A framework for the 3-d cloud monitoring based on data stream generation and analysis,” in *Integrated Network Management (IM), 2015 IFIP/IEEE International Symposium on*, pp. 62–70, IEEE, 2015.
- [8] J. Cito, P. Leitner, T. Fritz, and H. C. Gall, “The making of cloud applications an empirical study on software development for the cloud,” *arXiv preprint arXiv:1409.6502*, 2014.
- [9] R. Kohavi, R. M. Henne, and D. Sommerfield, “Practical guide to controlled experiments on the web: listen to your customers not to the hippo,” in *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 959–967, ACM, 2007.
- [10] “cloudwatchdev.” <http://awsdocs.s3.amazonaws.com/AmazonCloudWatch/latest/acw-dg.pdf/>.
- [11] “cloudwatch.” <https://aws.amazon.com/cloudwatch/>.
- [12] “azurewatch.” <http://www.paraleap.com/azurewatch>.
- [13] “newrelic.” <https://newrelic.com/>.
- [14] “aas.” <https://aws.amazon.com/autoscaling/>.
- [15] “nimsoft.” <http://www.nimsoft.com/solutions/nimsoft-monitor/cloud>.

-
- [16] “cloudyn.” <http://www.cloudyn.com/>.
- [17] “uptime.” <http://www.uptimesoftware.com/cloud-monitoring.php>.
- [18] “nagios.” <http://nagios.sourceforge.net/docs/nagioscore-3-en.pdf>.
- [19] “cloudsleuth.” <https://cloudsleuth.net/>.
- [20] “nimbus.” <http://www.nimbusproject.org/>.
- [21] “cloudstone.” <http://radlab.cs.berkeley.edu/wiki/Projects/Cloudstone>.
- [22] “groundwork.” <http://www.gwos.com/features/>.
- [23] “boundary.” <https://boundary.com/>.
- [24] “logicmonitor.” <http://www.logicmonitor.com/monitoring/storage/netapp-filers/>.
- [25] “cloudfloor.” <http://cloudfloor.com/>.
- [26] “cloudkick.” <http://www.cloudkick.com/home>.
- [27] “cloudclimate.” <http://www.cloudclimate.com>.
- [28] “monitis.” <http://portal.monitis.com/>.
- [29] “cloudharmony.” <http://cloudharmony.com/>.
- [30] “openstack.” <http://docs.openstack.org/diablo/openstack-compute/admin/oscompute-adminguide-trunk.pdf>.
- [31] C. C. Marquezan, D. Bruneo, F. Longo, F. Wessling, A. Metzger, and A. Puliafito, “3-d cloud monitoring: Enabling effective cloud infrastructure and application management,” in *Network and Service Management (CNSM), 2014 10th International Conference on*, pp. 55–63, IEEE, 2014.
- [32] C. Bunch, V. Arora, N. Chohan, C. Krintz, S. Hegde, and A. Srivastava, “A pluggable autoscaling service for open cloud paas systems,” in *Proceedings of the 2012 IEEE/ACM Fifth International Conference on Utility and Cloud Computing*, pp. 191–194, IEEE Computer Society, 2012.
- [33] A. Biswas, S. Majumdar, B. Nandy, and A. El-Haraki, “Predictive auto-scaling techniques for clouds subjected to requests with service level agreements,” in *Services (SERVICES), 2015 IEEE World Congress on*, pp. 311–318, IEEE, 2015.
- [34] “Scryer1.” <http://techblog.netflix.com/2013/11/scryer-netflixs-predictive-auto-scal.html>.

-
- [35] “Scryer2.” <http://techblog.netflix.com/2013/12/scryer-netflixs-predictive-auto-scal.html>.
- [36] “aws.” <https://aws.amazon.com/>.
- [37] L. R. Moore, K. Bean, and T. Ellahi, “A coordinated reactive and predictive approach to cloud elasticity,” 2013.
- [38] T. Lorido-Botran, J. Miguel-Alonso, and J. A. Lozano, “A review of auto-scaling techniques for elastic applications in cloud environments,” *Journal of Grid Computing*, vol. 12, no. 4, pp. 559–592, 2014.
- [39] D. Rowell, “State-space representation of lti systems,” URL: <http://web.mit.edu/2.14/www/Handouts/StateSpace.pdf>, 2002.
- [40] L. Ljung, *System identification*. Springer, 1998.
- [41] “arma.” https://en.wikipedia.org/wiki/Autoregressive%E2%80%9393moving-average_model.
- [42] “arx.” <http://zone.ni.com/reference/en-XX/help/372458D-01/lvsysidconcepts/modeldefinitionsarx/>.
- [43] “armax.” <http://zone.ni.com/reference/en-XX/help/372458C-01/lvsysidconcepts/modeldefinitionsarmax/>.
- [44] “oe.” <http://zone.ni.com/reference/en-XX/help/372458C-01/lvsysidconcepts/modeldefinitionsoe/>.
- [45] “bj.” <http://zone.ni.com/reference/en-XX/help/372458D-01/lvsysidconcepts/modeldefinitionsbj/>.
- [46] “Jpetstore.” <http://sourceforge.net/projects/ibatisjpetstore>.
- [47] “Agilefant.” <http://agilefant.com/>.
- [48] J. Cito, P. Leitner, H. C. Gall, A. Dadashi, A. Keller, and A. Roth, “Runtime metric meets developer-building better cloud applications using feedback,” *PeerJ PrePrints*, vol. 3, p. e1214, 2015.