



University of Science and Technology in Krakow

**Faculty of Computer Science, Electronics and Telecommunications
Department of Computer Science**

Autonomic cloud computing platform for scaling users' application

**DARIUSZ CHRZAŚCIK
RADOSŁAW MORYTKO**

Supervised by **MARCIN JARZĄB**
Assistant Professor of Computer Science

Krakow 2013

Oświadczamy, świadomi odpowiedzialności karnej za poświadczenie nieprawdy, że niniejszą pracę dyplomową wykonaliśmy osobiście i samodzielnie (w zakresie wyszczególnionym we wstępie) i że nie korzystaliśmy ze źródeł innych niż wymienione w pracy.

.....

PODPIS



Akademia Górniczo-Hutnicza im. Stanisława Staszica w Krakowie

**Wydział Informatyki, Elektroniki i Telekomunikacji
Katedra Informatyki**

Framework dla autoskalowalnego środowiska wykonawczego dla aplikacji użytkowników w konfiguracji hybrydowej chmury obliczeniowej

**DARIUSZ CHRZĄŚCIK
RADOSŁAW MORYTKO**

Promotor: dr inż. Marcin Jarzab

Krakow 2013

I would love to express my best wishes

..

Abstract

Cloud computing has become an attractive model for provisioning on demand computing resources as services to end-users. It is based on the assumption that almost anything can be viewed as a service, starting from applications delivered over Internet, through hardware in the data centers and ending on computing power. That model appears to be so attractive as from the user point of view the offered resources are infinite, transparent, robust and ready to consume at any time.

Most of the times end-users do not know in advance what the demand for the service is. This creates a requirement in which their systems are auto-scalable, i.e. they support sudden spikes in demand followed by underutilization at other times. An architecture of a cloud computing system, that meets these requirements, has a characteristic of a multi-hierarchical autonomic system, where different layers corresponds to different levels where cloud operates: starting from an application layer, through the application platform, infrastructure to end on a cloud instance.

Problem complexity raises challenges in a variety of aspects, especially in terms of providing cooperation and mutual sharing of resources that may belong to different kind of cloud providers. Therefore, an architecture that enables seamless cooperation among cloud providers and takes into account various QoS requirements of end-users be developed is absolutely vital. The InterCloud architecture is one of the first attempts that had been made in this direction. Having characteristic of an application platform in mind, we propose a variation of that architecture that supports cooperation among these application platform and fulfils needs of a decentralized environment.

The main contributions of this thesis are as follows: a) proposition of an architecture that enables aforementioned scenarios, b) implementation of that architecture c) simulation and laboratory tests.

Contents

1. Introduction	7
1.1. Motivation	7
1.1.1. Business potential	7
1.2. Contributions	8
1.3. Impact	9
1.4. Thesis structure	9
2. Scaling applications	10
2.1. Introduction	10
2.2. Horizontal scaling	12
2.2.1. Load-balancing algorithms	12
2.2.2. Load-balancing scalability	13
2.2.3. Load-balancer comparison	13
2.3. Vertical scaling	15
2.3.1. Virtual machine resizing	15
2.3.2. Virtual machine replacement	15
2.4. Providers comparison	16
3. Platform adaptivity	17
3.1. Introduction	17
3.2. Policies	18
3.3. Data analysis	18
3.4. Triggered actions	19
3.5. Providers comparison	19
4. Interoperability of clouds	21
4.1. Introduction	21
4.2. Hybrid cloud	21
4.2.1. Deployment models	21
4.2.2. Current usage and trends	22
4.3. Federation of clouds – InterCloud	23
4.3.1. Usage in industry	24
5. Design of Cloud-SAP	25
5.1. Requirements	25
5.2. High-level design	25

5.2.1. Overview	25
5.2.2. Auto-scaling subsystem	26
5.2.3. Cloud federation.....	29
5.3. Auto-scaling module	29
5.3.1. Application platform layer.....	29
5.3.2. Container layer.....	29
5.3.3. Stack layer.....	30
5.3.4. Cloud layer.....	30
5.4. Cloud federation	30
5.5. Solution discussion	30
6. Implementation	31
6.1. Requirements	31
6.1.1. Functional	31
6.1.2. Non-functional	32
7. Evaluation.....	33
7.1. Introduction	33
7.2. Cost of service deployment	33
7.3. Auto-scaling – single-provider based	36
7.4. Auto-scaling – multiple-provider based	43
7.5. Deployment time – solution comparison.....	49
8. Summary	53
A. Code listings.....	54
A.1. Service specifications	54
A.2. Scaling policies.....	56

1. Introduction

1.1. Motivation

One of the keys factors that has driven transformation of computing industry in the last years is the perception of computing utilities as an ordinary property[24], which can be easily accessed and adjusted to a specific needs. That point of view resulted in profusion of different services, often collectively referred as a cloud computing [36]. Similarly to services common to traditional markets, customers expect them to be accessible on demand and in easy manner, while paying only for the consumed goods. Furthermore, customers are interested in a given service provider only when it is eligible to guarantee appropriate quality of service.

The particular service providers that are addressed by this paper are the ones that supply users with an application execution platform, what is widely known as providing Platform-as-a-Service. In that case, a customer is an entity that has developed application and is eager to deploy it on an application platform that is able to fulfil his specific requirements, both in terms of quality and cost.

Having customer requirements in mind, it is crucial that service provider is able to adapt itself to meet them. For example, such adaptation can be triggered by a sudden spike in resource demand and may result in provisioning additional application platforms. However, due to the complexity of a system under consideration, there are different levels where adaptation is possible:

- user application
- application platform
- infrastructure

What is more, the fact that single service provider is constrained by his finite amount of resources poses a risk that it may not be able to serve customer all the time. Consequently, it is expected that adaptation at a service provider level is also possible, i.e. provider can offload some traffic to a different provider, as long as it satisfies a customer.

While autonomic computing has a long history [37], it has not been directly applied to a multi-layered problem that exists in a cloud computing environment. Especially, the research area at the last layer, which sizes across different service providers, is new. Although, architecture known as InterCloud [21] investigates problem of co-operation and negotiation at cloud level, it neither has been implemented nor presented in context of autonomic system.

1.1.1. Business potential

The rapid growth of interest in cloud computing in recent years resulted in huge sums of money being invested in the field. Figure 1.1 shows the size of the public cloud services market in 2012 and the forecast of its nearly two times growth in 2016. This data suggests that the subject is attractive for IT industry from the economic point of view. However, higher amounts of money spent on cloud services involve higher expectations of theirs quality

from customers. Although the most significant players in cloud computing have been in the field for quite a long time, it is still possible to outline some deficiencies their products have. Additionally, lack of common standards hinder cooperation among different cloud providers. For example it is nearly impossible to create an autoscaling cloud federation with Amazon Web Services (current leader in providing cloud services[32]) and another provider. What is more, Amazon AWS users cannot use more lightweight virtualization methods, such as linux containers. Nevertheless, when compared to other companies especially in terms of autoscaling capabilities, Amazon really shines. OpenShift, RedHat PaaS solution, ensures application scaling but with very limited possibilities of customisation of the process – the user can only choose if their application should scale and the whole algorithm is solely based on the number of concurrent requests to the application. Users of Heroku, another PaaS solution, have no automation tool that would control the number of instances (*dynos* in Heroku nomenclature) their application is running on – they can change it manually.

The proposed solution in this dissertation tries to deal with the aforementioned providers problems by outlining an example architecture that enables seamless cooperation among cloud providers and provides auto-scaling capabilities.

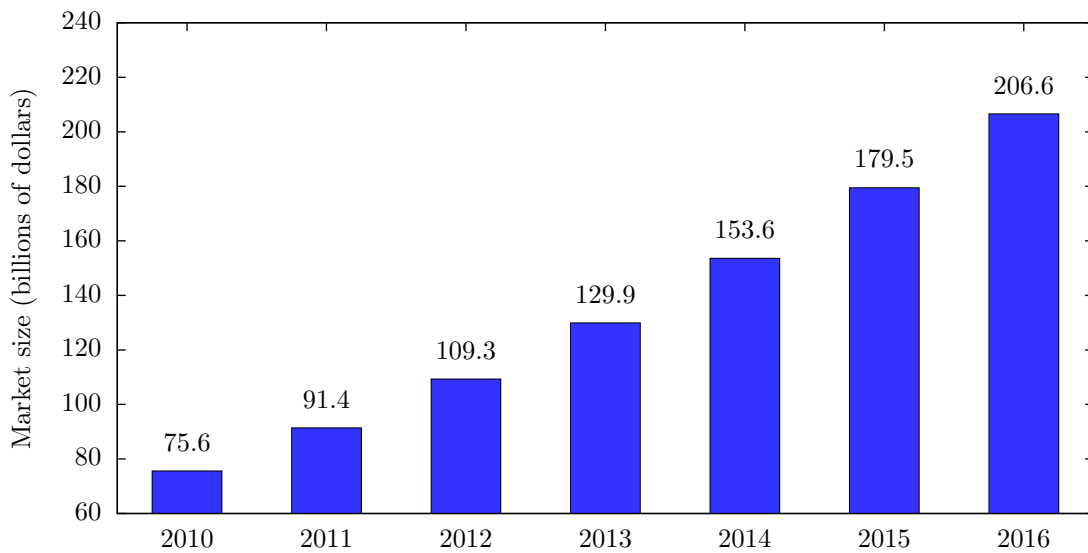


Figure 1.1: Public Cloud Services Market Size, 2010-2016 (forecast). Source: *Gartner*, 08/2012

1.2. Contributions

The main contributions of this dissertation are as follows:

- A proposal of an architecture of federated cloud computing environment, which is based on and can be viewed as a simplified version of *InterCloud*
- The notion of considering each service model as an autonomic system
- The implementation of the proposed architecture using OpenNebula technology stack

1.3. Impact

We hope that the concept of representing each level of an autoscaling subsystem as an autonomic one can be thought-provoking for cloud computing scientists. What is more, we believe that our successful attempt to implement a simplified variant of an InterCloud architecture will cause its gain in interest and popularity. Finally, we consider the ideas contained in this work be beneficial to the OpenNebula ecosystem as they provide insights into the ways Quality of Service can be ensured:

- implementing autoscaling capabilities
- designing *cloud infrastructure* in accordance with InterCloud architecture

1.4. Thesis structure

2. Scaling applications

This chapter is devoted to the concept of scaling users' application from the perspective of a cloud platform provider. To achieve that, it presents attainments of research groups working in that area as well as it considers mechanisms used in products currently available on the market.

2.1. Introduction

The reason why scaling application lies in our area of interest is the fact that it is widely accepted measure for improving application performance, consequently increasing offered Quality-of-Service. Enriching system with capability to scale entails avoiding additional costs that are related to coping with excessive traffic. In some cases, these costs may be caused by not handling extra traffic at all and may involve aspects such as: increased response time, processing overhead, space, memory, or money [18].

While scalability is a widely used term, it still lacks a clear and concise definition. Over the time, there were a few attempts to define it, yet not all of them were claimed as successful [28] [26]. Hence, it is necessary to clarify this term before going into further discussion. Instinctively, scalability is perceived as ability of a system to accommodate an increasing number of elements or objects to process. In particular, we can point out different types of scalability that are affected by increased number of requests: [18]:

- *load scalability* - ability to work without delays and unproductive resource consumption at light, moderate, or heavy loads while making good use of available resources. Factors that may hinder load scalability include: scheduling shared resource, self-expansion, inadequate exploitation of parallelism
- *space scalability* - memory requirements do not grow to intolerable levels as the number of items system supports increases
- *space-time scalability* - system continues to function gracefully as the number of objects it encompasses increases by orders of magnitude
- *structural scalability* - implementation or standards do not impede the growth of the number of objects system encompasses

Although, all of the aforementioned aspects are vital for any application, our work focuses solely on the first type of scalability. The reasoning behind this statement is that, while all of these scalability types lies in direct responsibility of an application developer, the load scalability can be additionally improved by adding additional resources to a system. This brings us to a question what kind of resources are used by an application or more appropriately in context of this dissertation: *what kind of resources can we add to improve application performance?* Required resources varies from an application to an application. However, among the most common ones we can distinguish:

- CPU

- memory
- storage
- network bandwidth

It is commonly agreed that there are two main possible ways the resource can be added:

- *horizontal scaling (scaling out)* - adding more nodes to a system, such as servers in a context of distributed application
- *vertical scaling (scaling up)* - increasing capacity of a single node in a system, i.e. adding additional memory, CPU, storage, etc.

What makes scaling application particularly interesting are the benefits offered by a cloud computing, especially the illusion of a virtually infinite computing infrastructure [40]. Making use of virtualization technologies, which often underpins cloud computing platform, allows for resource manipulation in a dynamic, on-demand manner. Although, cloud computing offers additional scaling capabilities, it increases solution complexity since they operate in different layers: server, platform, network as stated in [40]. However, since platform containers are often represented either as virtual machines or another isolated environment (e.g. OpenShift leverages SELinux and cgroups) they are similar in nature to server scaling and supports both scaling up and out. Therefore, the remaining of this chapter is focused solely on server scaling, omitting network scaling as it lays outside of scope of this dissertation.

Having that said, common sense dictates that adding resources is only a part of the success - it should be accompanied by tuning application platform configuration. For example, adding supplementary CPUs without increasing thread pool size that handle requests makes a little sense. Similarly, in context of a Java application, we have to increase heap size, to make a good use of extra memory. While importance of application tuning cannot be underestimated, its detailed analysis lies outside of the scope of this dissertation. Figure 2.1 presents different scalability layers and actions that can be taken at each level to improve application performance.

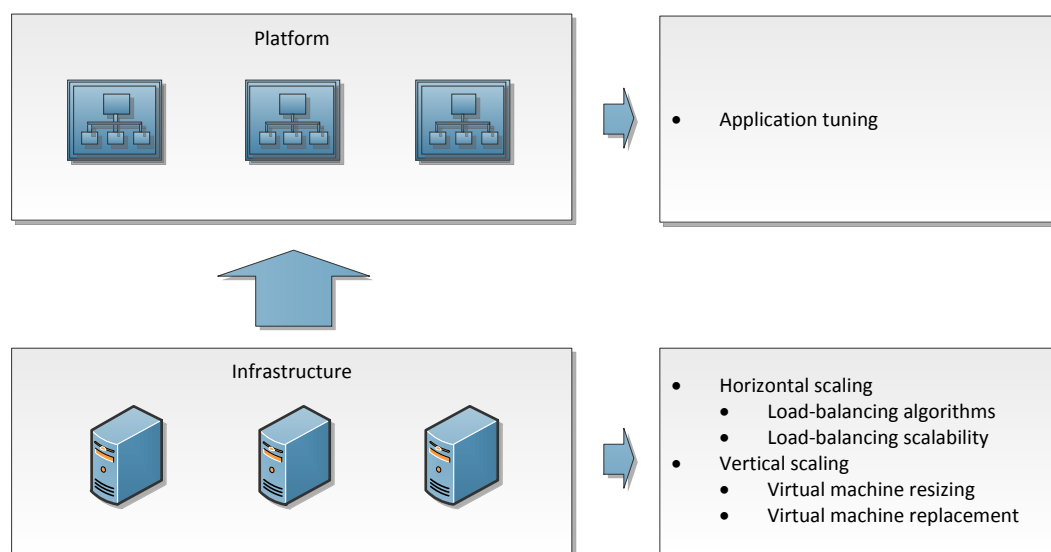


Figure 2.1: Scalability layers

With all that said, there is no silver bullet - not matter what underlying mechanism platform provider decides to use, the application developer is still responsible for creating an application with scaling in-mind. This statement has been already proven in 1967 by Amdahl law, which in short states that sequential component of a parallel algorithm impacts efficiency for a sufficiently large number processors [27] as shown in Figure 2.2. In other words, adding supplementary resources to a poorly written application (i.e. having a lot of sequential or synchronized components) can be beneficial only to a certain degree.

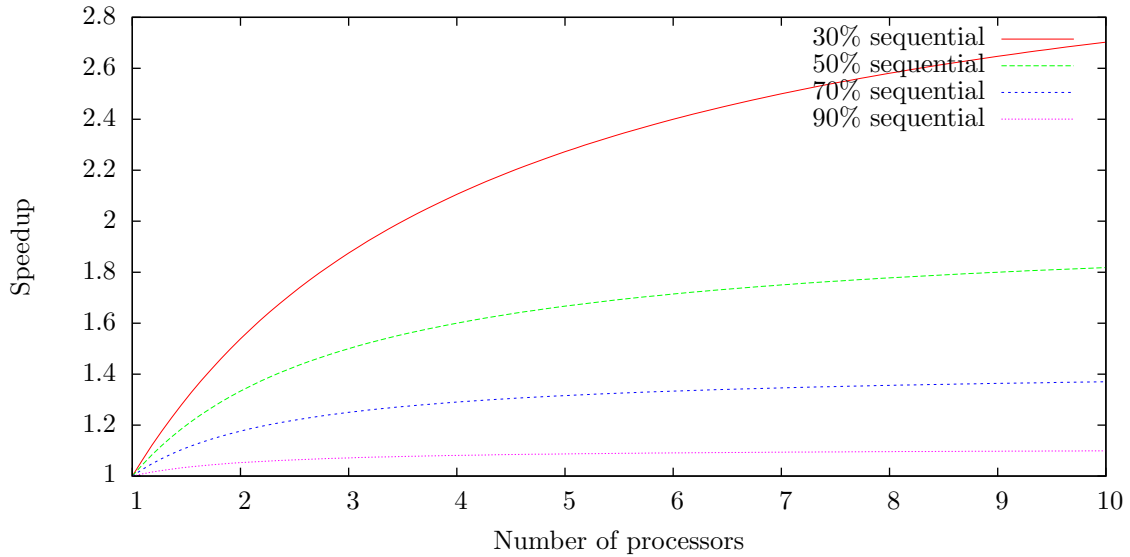


Figure 2.2: Amdahl's law

The rest of this chapter elaborates in detail about horizontal and vertical scaling taking into account mechanisms used in Platform-as-a-Service solutions that are available on the market.

2.2. Horizontal scaling

As outlined in previous section, horizontal scaling is about adding supplementary nodes to a system. As it is common to cloud computing, nodes are represented as virtual machines and this assumption is used in further discussion. Consequently, adding server comes down to cloning a new virtual machine from a template and possibly installing additional software and reconfiguring it later. While mechanism of creating new virtual machine from a template is offered literally in every IaaS platform currently available (OpenStack [15], OpenNebula [12], CloudStack [5] or Eucalyptus [7] to name a few) and is similar in manner, the underlying hardware and virtualization mechanism determines how fast provisioning is done.

Provisioning new server is only a first step in scaling an application, it is required to configure load balancing mechanism to make use of additional node. The two important aspects that have to be consider are: load-balancing algorithms and scalability.

2.2.1. Load-balancing algorithms

Generally, there are two types of load-balancers: hardware and software based. Due to the dynamic nature of system under consideration, we focus only on the latter as it offers a greater deal of flexibility. Among the most common algorithms we can distinguish [10]:

- *round-robin scheduling* - request are sent to successive nodes, according to their weights. This algorithm is fairest when the server's processing time remains equally distributed [10]
- *least connection* - the server with the lowest number of connections receives the connection
- *source routing* - source IP address is hashed, the same client IP address always reaches the same server
- *URI hashing* - URI that designates resource is hashed and divided by the total weight of the running servers. Such hash designates which server that receives the request. In practice, this algorithm is commonly used with proxy caches and anti-virus proxies in order to maximize the cache hit rate.
- *request counting algorithm* - load is distributed the requests among the various workers, ensuring that each gets their configured share of the number of requests
- *weighted traffic counting algorithm* - variation of above-mentioned algorithm with a difference that it is focused on bytes rather than number of request
- *pending request counting algorithm* - scheduler keeps track of how many requests each worker is assigned at present. A new request is automatically assigned to the worker with the lowest number of active requests

Situation gets further complicated when considering real-world web application that sends user information using cookies, what imposes requirement on load-balancer for session stickiness [39].

2.2.2. Load-balancing scalability

Although, it may seem that balancing workloads eliminates problem of a single point of failure (SPOF) among different servers, it is in fact shifted to load-balancing layer. In other words, load-balancer becomes a new SPOF. Therefore, in cases where high availability is required, multi-tiered load balancing architecture should be considered. This, however, seems not to be a case among IaaS or PaaS providers - none of them unequivocally specifies whether they provide redundancy at load-balancer level.

2.2.3. Load-balancer comparison

While there are many load-balancers available on the market, following are credited to be most popular:

- **HAProxy** [9] - load-balancer initially written by Willy Tarreau. Noticeably, it's used by OpenShift [13] to distribute load among gears [14]
- **BIG-IP Local Traffic Manager (LTM)** - solution offered by F5 [8]. Although LTM is a hardware solution, omitted in this section, it also has also its virtualized counterpart.
- **Apache HTTPD** [1] - popular HTTP server. When enhanced with additional modules, it can behave like a proxy or load-balancer. Over the time, there were several attempts to develop such modules: `mod_jk` [2], `mod_proxy_balancer` [3], to name a few. While the former is purely AJP13 oriented, the latter supports different protocols: HTTP, FTP and AJP13. As a consequence, only `mod_proxy_balancer` was taken into account during comparison.

Table 2.1 presents they key performance features and algorithm used to schedule requests.

	Performance features	Scheduling algorithms
HAProxy	<ul style="list-style-type: none"> – a single-process, event-driven model reduces the cost of context switch and the memory usage – O(1) event checker – single-buffering without copying data between reads and writes – zero-copy forwarding – optimized HTTP header analysis: headers are parsed and interpreted on the fly 	<ul style="list-style-type: none"> – round-robin scheduling – least connection – source routing – URI hashing
BIG-IP Local Traffic Manager	<ul style="list-style-type: none"> – managing at application services level rather than at individual devices and objects – scripting language that allows administrator to intercept, inspect, transform, and direct application traffic – built-in firewall protection, application security, and access control – real-time protocol and traffic management decisions 	
Apache HTTPD	<ul style="list-style-type: none"> – support for session stickiness by using cookies and URL encoding. This approach [3] avoids unequal load distribution if clients are hidden behind proxies and stickiness errors when a client uses a dynamic IP address that changes during a session 	<ul style="list-style-type: none"> – request counting algorithm – weighted traffic counting algorithm – pending request counting algorithm

Table 2.1: Comparison of load balancers

2.3. Vertical scaling

Essentially, vertical scaling is concentrated upon increasing capacity of single node. Again, when considering technical advancements that comes with cloud computing and virtualization, we can differ two categories of scaling: virtual machine resizing and virtual machines replacement. This distinction is dictated by limitation hypervisors - not all of them are able to resize virtual machine without shutting it down.

2.3.1. Virtual machine resizing

	Memory	CPU	Disk
KVM 1.2.0		– dynamic pinning CPU to a specific virtual machine (depending on underlying hardware)	– adding a disk to a LVM group
Xen 4.3	– changing the amount of host physical memory assigned to virtual machine without rebooting it – start additional virtual machines on a host whose physical memory is currently full, by automatically reducing the memory allocations of existing virtual machines in order to make space	– dynamic pinning CPU to a specific virtual machine (depending on underlying hardware)	– dynamic block attaching, adding a disk to a LVM group
VMware ESX 5.1	– hot-plugging memory, ex. using VMware vSphere	– hot-plugging CPU, ex. using VMware vSphere	– adding additional disks to existing virtual machine
OpenVZ (kernel: 042)	– configurable via user beancounters	– configurable via user beancounters	– configurable via user beancounters

Table 2.2: Comparison of hypervisors resizing capabilities

2.3.2. Virtual machine replacement

As it was highlighted in previous section, reasoning behind virtual machine replacement is that, in case when dynamic resizing is not possible, a new virtual machine with a desired configuration can be provisioned and replace the old one. Since this is a basic operation, all above-mentioned hypervisors supports this scenario as long as required resources are available.

2.4. Providers comparison

Table 2.3 presents a summary of cloud providers auto-scaling capabilities. Interestingly, all of them are focused solely on horizontal scaling, ignoring advantages offered by a fine-grained approach to scaling that leverage scaling up and application tuning.

	Horizontal scaling	Vertical scaling	Application tuning
Infrastructure provider			
Carina	✓	×	×
OneFlow	✓	×	×
AWS EC2	✓	×	×
Platform provider			
CloudFoundry	×	×	×
OpenShift	✓	×	×
AppEngine	✓	×	×
Azure	✓	×	×
Heroku	×	×	×

Table 2.3: Comparison of cloud providers scaling capabilities

3. Platform adaptivity

This chapter introduces concepts and mechanisms that enhance a platform with adaptivity capabilities, which are often achieved by fusion of rules, policies and scaling techniques.

3.1. Introduction

In short, platform adaptivity adds a auto-scaling features to a solution offered by Platform-as-a-Service provider. Key concept is to have a Elasticity Controller which gathers probes from virtual machines and uses that knowledge to execute appropriate action on cloud instance, indirectly modifying consecutive probes [40]. This concept illustrates diagram 3.1.

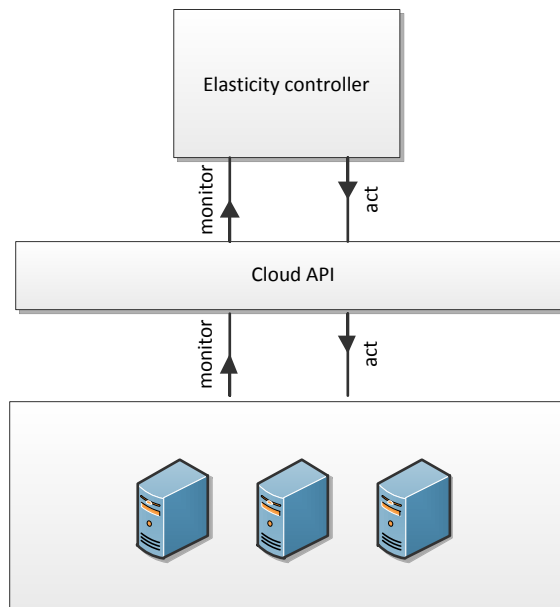


Figure 3.1: Elasticity controller

The remaining of this chapter describes crucial elements that compose elasticity controller: policies, data analysis, triggered actions and presents a comparison of cloud providers.

3.2. Policies

While auto-scaling is offered by a vast amount of cloud providers (e.g AWS, OpenShift, OpenNebula) it often lacks a sophisticated mechanisms allowing for specific scaling policies, being limited to only one predefined rule as it is in case of OpenShift for example.

Policy denotes a condition which, when satisfied, triggers an action that is supposed to harness cloud instance in a way that future evaluations of condition will be unsuccessful. Typically condition itself is accompanied by a minimal and maximal number of node instances, allowing for ensuring minimal QoS and controlling maximal costs. Currently, industry leaders supports two main kind of policies [4]:

- *expression based* - allows to define how you to scale application in response to changing conditions, which include factors such as memory, CPU usage, cost or some indirect, calculated metrics
- *scheduled* - allows to scale an application in response to predictable load changes. For example, traffic increases during the weekends and decreases on working days. Hence, that predictable traffic patterns is used to scale application based on current time.

Technically, policies are expressed in some human-readable format such as JSON, XML as it is in case of AWS EC2 or custom expression used for example by Carina. Appendix A.2 presents example configuration used by AWS E2 Auto-Scaling.

3.3. Data analysis

Having policies defined, their conditions are evaluated against data acquired from sensors. In a simplest case this evaluation can be based on a Threshold Model [33], which defines a valid range. In cases when given metric violates that condition (i.e. value is either smaller than minimal or bigger than maximal acceptable) corresponding resource is properly adjusted - figure 3.2 illustrates that idea.

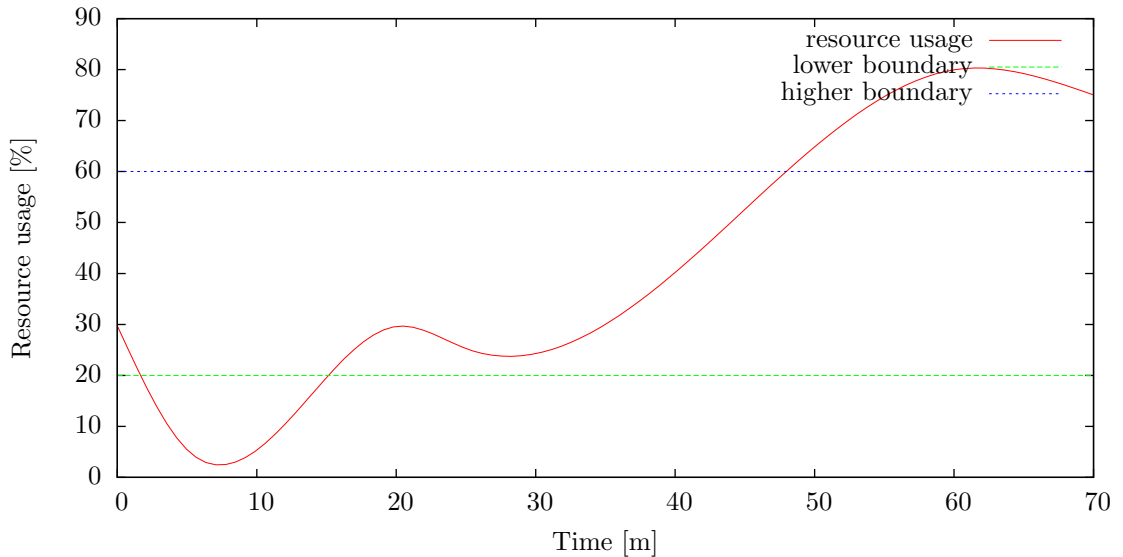


Figure 3.2: Threshold model

While trivial in its form, cases of AWS, OpenShift, Carina, OneCloud proves it is useful in a real-world scenarios. Having that said, more sophisticated algorithms also exists:

- *integer programming* - auto scaling is reduced to server integer programming problems, which aims to minimize the cost or maximize the computing power with either computing power constraints or budget constraints [34]
- *burst based padding* - employs a signal processing technique based on fast Fourier transform, burst pattern is extracted and used to calculate a padding value. Coefficients that represents the amplitude of each frequency component are used to calculate burst density. Depending of that value (i.e. is higher than 50%) appropriate percentile of the burst values are used [38]
- *remedial padding* - padding errors are being recorded and used in successive padding evaluations. In other words, let e_1, e_2, \dots, e_k denote the recent prediction errors, next the weighted moving average is calculated. Actual applied padding is either padding itself or weighted average, depending which one is greater [38]
- *linearised dynamic control* - linearised correction model is based on control equations [20]:

$$x(k+1) = Ax(k) + Bu(k) \quad (3.1)$$

$$y(k) = Cx(k) + Du(k) + Ez(k) \quad (3.2)$$

where x denotes the state variable vector and coefficient matrices: $A B C D E$ are fitted to historical data as a regression model - [25]

- *markov decision process model* - computes optimal reaction to state changes by using observation on the system with assumptions about the rate of changes expected in the future [19]

3.4. Triggered actions

Actions that are being triggered by a elasticity controller are focused on application scaling. Previous chapter described that problem in detail.

3.5. Providers comparison

Table 3.1 summaries chapter with approaches to adaptivity taken by different cloud providers. Section 'triggered action' is omitted for brevity - it was described extensively in previous chapter.

	Policies	Data analysis
Infrastructure provider		
Carina	<ul style="list-style-type: none"> – time frame based – expression based (only for CPU) 	– threshold model that takes into account minimal and maximal permitted instances of an application as well as application priority
OneFlow 4.2	<ul style="list-style-type: none"> – time frame based with customizable padding – expression based build on custom language, where all vm's metrics are supported – customizable adjustment padding, cooldown time 	– threshold model
AWS EC2	<ul style="list-style-type: none"> – time frame based – expression based, where expressions corresponds to a AutoScalingGroup – actions are triggered by a CloudWatch alarms – customizable adjustments paddings, types, cooldown time 	– threshold model, takes into account minimal and maximal permitted instances of an application as well as application priority
Platform provider		
CloudFoundry	×	×
OpenShift	– single built-in policy	– single built-in threshold model that scales an application when CPU load is greater than 50% for a given period
AppEngine	<ul style="list-style-type: none"> – built-in policy based on request queue length – adjustable minimal, maximal number of application instances, pending latency 	– queue-based, new instance is provisioned if queue length got too long
Azure	<ul style="list-style-type: none"> – time frame based – expression based, where expression can involve either CPU usage or Queue length – customizable adjustments paddings, types, cooldown time 	– threshold model, that takes into account minimal and maximal allowed instances
Heroku	×	×

Table 3.1: Comparison of cloud providers approach to adaptivity

4. Interoperability of clouds

This chapter introduces the notion of a hybrid cloud and explains its role in IT industry. On top of this deployment model, the concept of InterCloud is presented and elaborated with the emphasis on its application in ensuring scalability of users' services.

4.1. Introduction

From the perspective of a user of PaaS services it is vital that they are able to deploy seamlessly their applications using libraries, tools and services supported by the cloud provider[36]. Judging by such factors as the popularity of Heroku – currently one of the most popular PaaS providers which does not offer more advanced features which would enable management of the infrastructure underpinning the deployment platform, the fact that Microsoft added auto-scaling to its Azure platform as late as in June 2013, it is perfectly possible most PaaS users are satisfied with the current offers of their providers and do need another, more sophisticated functionalities. However, there are more complex applications and systems whose requirements regarding technology stack, availability and scalability are considerably more demanding. For such services there ought to be designed slightly specialized features that would require cooperation among different cloud providers.

4.2. Hybrid cloud

One can imagine scenarios in which customers of cloud services know their applications are vulnerable to sudden variations in demand and their responsiveness must be kept at the same level all the time. In such cases, they want them to scale dynamically according to current load or other predefined or manually specified metrics. What is more, in order to ensure high availability of their services, customers do not want to confine themselves to only one provider – in the best scenario they want their applications (or their logical parts, such as persistence layer) to be spanned across different providers and be able to cooperate with one another at the same time. Additionally, due to privacy concerns of the sensible data, companies are reluctant to put it in the public cloud storage. All these factors lead to the concept of a *hybrid cloud*[36] – the case in which the cloud is a composition of two or more distinct infrastructures which are unique entities, but there are technological means that make it possible to port data and applications among them.

4.2.1. Deployment models

The informal introduction to the concept of a *hybrid cloud* in the previous section requires a strict definition, but it is virtually impossible without defining other deployment models:

- Private Cloud – The provisioned cloud infrastructure is used exclusively by a single organization (that may consist of many business units) and may be owned, managed and operated by the organization or a third party.
- Public Cloud – The provisioned cloud infrastructure is used by general public and may be owned, managed and operated by a business, academic or government organization or some combination of them. It exists on the premises of the cloud provider.
- Community Cloud – The cloud infrastructure is provisioned for exclusive use by a specific community of consumers from organizations that have shared concerns (e.g., mission, security requirements, policy, and compliance considerations). It may be owned, managed, and operated by one or more of the organizations in the community, a third party, or some combination of them, and it may exist on or off premises.

Having defined those models, we can see that *hybrid cloud* can be placed among them and be defined as a model in which the provisioned infrastructure is a composition of two or more other infrastructures - *private, community* or *public*.

4.2.2. Current usage and trends

Cloud – clients' view

Before digging into the details of current usage and popularity of the hybrid model, it is worth discussing the general attitude of clients towards cloud computing. As the recent survey [16] shows, the major factor that prevents companies from adopting cloud solutions is their concern over security – in 2012 as much as 52% responders considered it as a main concern with the regard to cloud in general. However, the tendency is that more and more enterprises do not find it a major issue as in 2012 the number declined to 46%. Complexity related to the management of cloud components, Vendor lock-in, interoperability and reliability were among the most frequent obstacles to adoption in 2013 for they constituted 46%, 35%, 27% and 22.3% of responders' votes respectively. Total results are shown in the figure 4.1.

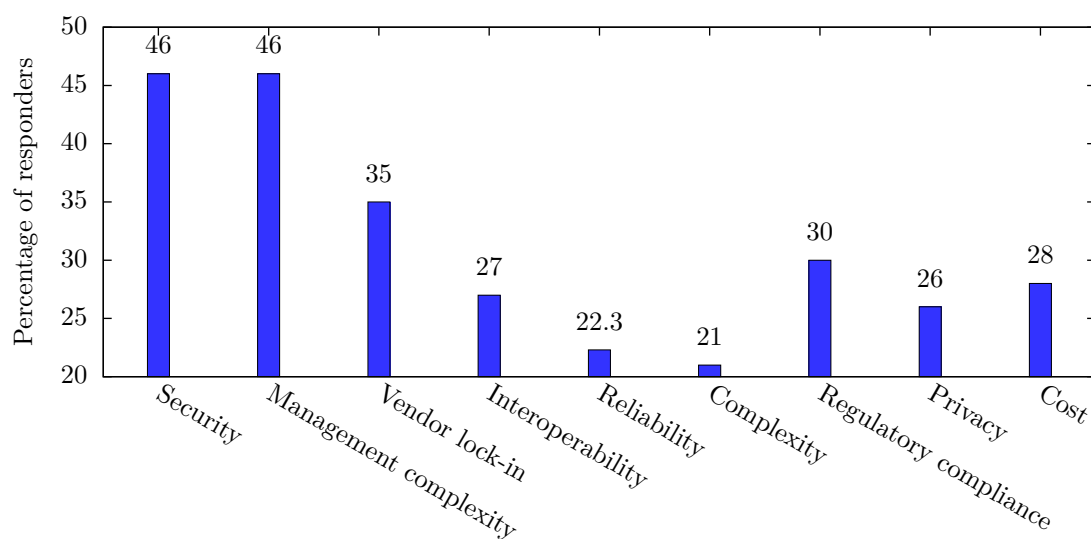


Figure 4.1: Major obstacles to cloud adoption in 2013 according to [16]

The same survey shows that the cloud adoption growth rate is high – 75 percent of responders stated usage of some sort of a cloud platform. This means 8 percent growth when compared to the results obtained in 2012. The expectations for the total worldwide addressable market for cloud computing are to reach \$158.8B by 2014 – an increase of 126.5 percent from 2011.

View on hybrid cloud

When it comes to the application of a hybrid model in industry, in most cases the definition introduced in the previous chapter now becomes a 'public-private' composition. And this is how the term should be understood when discussing the results of the surveys which aimed to provide insights onto the view on a hybrid cloud from the customers' perspective. The study [16] forecasts 16 percentage growth in the hybrid cloud adoption in 5 years, from 27 to 43 percent. At the same time, the usage of a public model will decline from 39 to 32 percent. The other survey, conducted by Rackspace [17], provides more detailed data about current usage and popularity of a hybrid model. The first interesting finding is that as much as 60% responders, which included 1300 companies in the UK and US, have moved or are planning to move certain applications either partially (41%) or completely (19%) off the public cloud because of its limitations or the potential benefits of other models, e.g. the hybrid one. The second one is about the pros of adopting the hybrid cloud – potential users find more control (59%) and better security (54%) top benefits of using this deployment model. The other most responded benefits are shown in the figure 4.2.

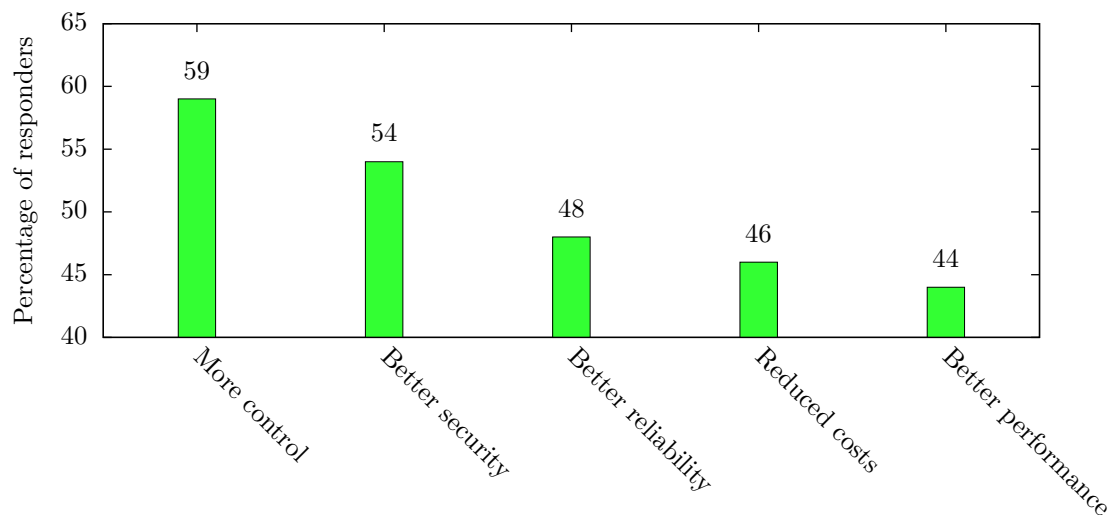


Figure 4.2: Top benefits of using the hybrid model according to [17]

4.3. Federation of clouds – InterCloud

As stated in the previous sections, one of the major obstacles that prevents consumers from adopting cloud solutions is reliability. It appears that these statements are not only imaginary worries of entrepreneurs, but real problems – there are cases, where some providers temporarily run short of capacity (e.g. because of provisioning too many virtual machines) in the face of high demand [31]. What is more, the more demanding clients require specific QoS to be satisfied by their providers as negotiated in Service Level Agreements. In order to meet these challenges there is a need of a completely new approach to the problem of effective management of resources. The new solution should take into account such factors as:

- users’ requests priority
- users’ QoS requirements (e.g. the deadline by which some jobs have to be executed)
- price that the clients pay for the usage of resources

Computer scientists in the field of cloud computing devised a model [23] in which resources are managed in a market-oriented fashion that enables dynamic regulation of the supply and demand of resources and promotes the mechanisms for their allocation that would take into account their priorities and levels of utilization. The extension of this model is a vision of creating the federated cloud computing environment, so called *InterCloud*, that “facilitates just-in-time, opportunistic, and scalable provisioning of application services, consistently achieving QoS targets under variable workload, resource and network conditions” [22]. The elements of the proposed architecture are as follows:

- Cloud Exchange – acts as a market maker for bringing together both producers and consumers of services. It allows Cloud Brokers and Cloud Coordinators to match consumers with the fitting offers from providers. Such a market is a step forward towards creating a dynamic infrastructure for trading based on Service Level Agreements.
- Cloud Coordinator – manages the instance of a cloud and its membership in the overall federation; provides an environment (programming, deployment) for applications
- Cloud Broker – acts on behalf of the client; communicates with the Cloud Exchange to find the best cloud instances for the application

4.3.1. Usage in industry

The depicted model has not yet been adopted in the industry, yet some simulations were carried out on a *CloudSim* platform and the obtained results showed that this concept has “immense potential” [22].

5. Design of Cloud-SAP

This chapter introduces the high-level design of Cloud-SAP, highlighting its core concepts and indicating possible implementation ideas.

5.1. Requirements

One can notice that elements that yields a solution for a problem stated in the first chapter, which is ensuring that users' application provide appropriate Quality-of-Service for its customers in a most-cost effective manner, were gradually introduced in previous chapters:

- *scalability* - ability to improve application performance by enriching resources
- *adaptivity* - ability to adapt (i.e. scale) appropriately to a current usage pattern
- *inter-cloud awareness* - ability to compose an application deployment using different cloud providers; cooperation with different cloud provider to supply application with extra resources while performing application scaling

Next section states the general overview of the proposed solution, while the consecutive sections details its elements and finally the last section summarises the design choices in a context of system requirements.

5.2. High-level design

5.2.1. Overview

Diagram 5.1 illustrates exemplary Cloud-SAP deployment, among depicted components we can distinguish two main parts of proposed solution: auto-scaling subsystem and inter-cloud broker. Auto-scaling subsystem is responsible for scaling users' application, taking into account different scalability perspectives, handled by a different layers:

- application platform layer: application platform tuning
- container layer: vertical scaling
- stack layer: horizontal scaling
- inter-cloud layer: scaling out across different cloud providers

The last mentioned layer - inter-cloud - cooperates with different cloud providers, hence, it requires additional work to be done. In fact, these extra tasks are delegated to a inter-cloud broker, responsible for a dialog in a cloud federation.

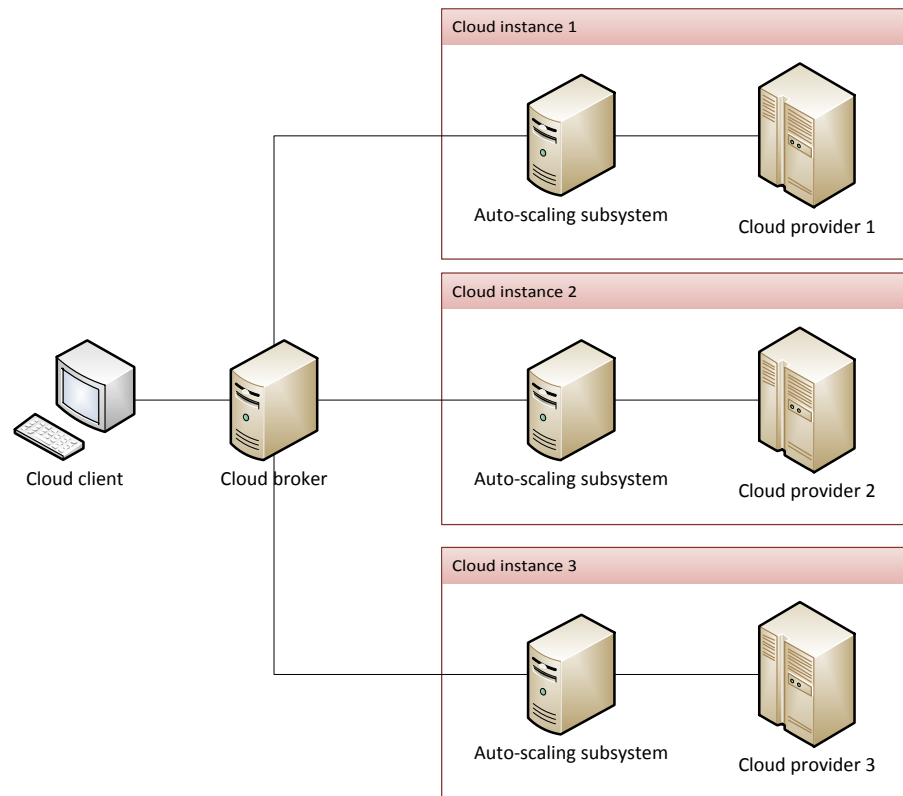


Figure 5.1: Cloud-SAP high level overview

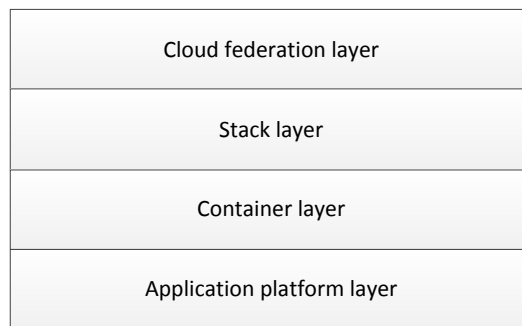


Figure 5.2: Layered structure of Cloud-SAP

5.2.2. Auto-scaling subsystem

Autonomic components

Each layer of the Cloud-SAP (shown in figure 5.2) must be characterised by an ability to adapt to a given application usage. One of the first models that ensured system adaptivity is an autonomic component, concept based on a feed-back loop, initially proposed by IBM [29]. Figure 5.3 depicts that architecture. This observation is a foundation of proposed architecture - adaptivity is achieved by enriching each layer with an elasticity controller that is in fact an autonomic component.

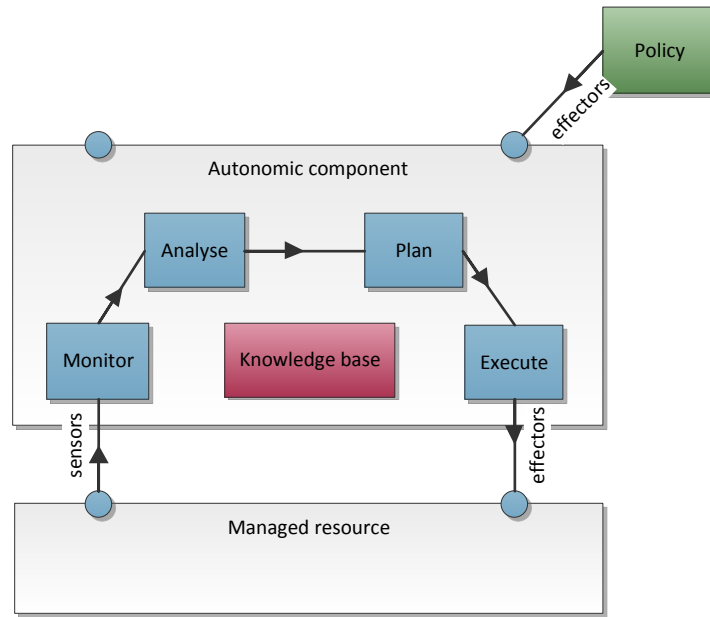


Figure 5.3: autonomic component

Due to the fact that designed platform operates on multiple layers, we can extend idea of an autonomic component by using concept of a multi-hierarchical autonomic system [33]. Figure 5.4 illustrates that hierarchy, where each level represents a different perspective on an application scaling. First three hierarchy levels (application tuning, container, stack) are centralised and controlled by a single elasticity controller at each level, while the last inter-cloud level is a decentralized one - each cloud instance is fully independent.

Considering hierarchy of our architecture, each level manages an underlying autonomic component, while being managed by an upper layer at the very same time.

Monitoring

Monitoring module of each elasticity controller aims to gather data from *Sensors* of a managed resource. While design of Cloud-SAP doesn't have any specific requirements regarding underlying monitoring mechanism, there are a few aspects that should be taken into consideration. Table 5.1 summarises main points.

Cloud-SAP is entirely focused on further data processing, consequently decision whether to store monitoring data or not is entirely implementation dependent.

Analysis

Analyse module applies a specified mathematical model to a gathered data and yields a conclusion on top of that. There is a vast number of models that can be used for that purpose, for example [33] enumerates:

- Queue-based Performance Models
- Dynamic models
- Monotonic static models
- Policy based models

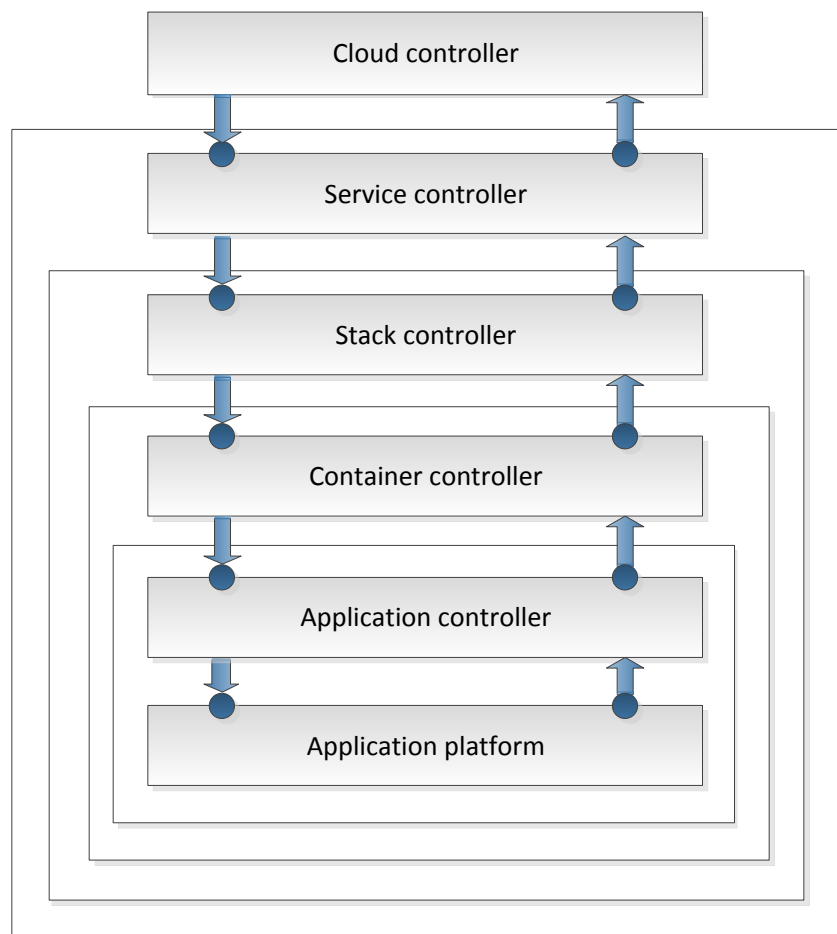


Figure 5.4: Cloud-SAP as an hierarchical autonomic system

It is vital that Cloud-SAP implementation use predictive approach [30] and track changes [41] in a system, giving a solid foundation for a knowledge base that can contribute to a better data analysis in future.

Planning

Planning module is responsible for scheduling actions that have to be taken to solve problems reported by an earlier analysis. This may involve:

- reserving some resources (for example to deploy a new virtual machine)
- handling situation where multiple entities compete about the resources (for example taking service's priority into account)

It is possible, however, that resolving issues is not possible at a current layer. In that case, the upper-layer is responsible for handling it.

Execution

It is expected that execution module manage resource by performing various operations on resource's effector. Each layer is characterised by different actions that can be taken, for example: application server reconfiguration,

	Advantages	Disadvantages
data exchange model		
push based	– notification driven: all changes are reflected immediately in the system	– requires additional agent at monitored component
pull based	– more calculations are being done at monitoring module side (doesn't involve additional load at managed component)	– some monitoring cycles can be redundant if scheduler is not properly tuned
data format		
binary	– space-effective	– unreadable for human
character-based	– human readable	– ineffective in terms of space
communication protocol		
vendor specific	– provides specific data to a given managed component (ie. virtual machine / container)	– may cause incompatibility problem while expanding cloud-sap system to a consecutive cloud instances
standard-based	– supports scaling across multiple cloud instances as long as they support standard	– gathered data may be inadequate to a specific needs

Table 5.1: Monitoring - summary of issues

vertical or horizontal scaling.

Enforcing actions on a managed component can be done using different methods, however, all of them requires an agent installed in a managed resource. Cloud-SAP does not enforce any specific technology, leaving low-level decision to an implementation.

Knowledge base

Knowledge base is represented as a set of rules, policies shared across autonomic components.

– TODO przykład z ibma

5.2.3. Cloud federation

Opis najwyzej warstwy, ze ma brokery, wykorzystuje kilka chmur.

5.3. Auto-scaling module

5.3.1. Application platform layer

5.3.2. Container layer

Among the currently available monitoring systems remarkable are Cloud Watch [6], OpenNebula (Information Manager / Virtual Machine Manager) [12] there are also open approaches such as Ganglia Monitoring System [35].

5.3.3. Stack layer

5.3.4. Cloud layer

5.4. Cloud federation

5.5. Solution discussion

6. Implementation

In this chapter we outline implementation details about each component of the proposed solution.

6.1. Requirements

6.1.1. Functional

One can notice that elements that yields a solution for a problem stated in the first chapter, which is ensuring that users' application provide appropriate Quality-of-Service for its customers, were introduced in previous chapters:

- scalability – ability to improve application performance by enriching
- adaptivity – ability to adapt (i.e. scale) appropriately to current usage pattern
- inter-cloud awareness – ability to cooperate with different cloud provider to supply application with extra resources

Having those in mind, we can make the list of functional requirements more formal:

1. The user of the platform is able to:
 - (a) deploy a service,
 - (b) cancel the service,
 - (c) check the status of the previously ordered-to-deploy service at any time. *Status* means a) whether or not the deployment succeeded, b) current uptime of the service, c) current cost
2. One of the elements of the platform is a client application that is used by the user of the platform to communicate with it,
3. During the deployment process, the platform takes as an input a description of the service (application) that consists of:
 - service name,
 - software stacks (e.g. *java*, *ruby*),
 - auto-scaling policies (per each stack) which define i) minimal and maximal number of VMs that are needed for the stack, ii) name of the policy (algorithm) which is used for scaling, iii) parameters of the policy
4. Deployment of a service is done in a way which minimizes the cost from the client's perspective with ensuring Quality-of-Service requirements at the same time,

5. It is assumed that the application which is going to be deployed is properly and fully tuned so that it is not possible to improve its performance by changing its or any of its components configuration(s),
6. The platform monitors the state of the deployed services and based on the results of this process takes appropriate steps in order to meet the auto-scaling requirements. These include a) altering VM's parameters and configuration, b) vertical scaling, c) horizontal scaling, d) scaling stacks among different cloud providers

TODO Alternative scenario – the client has a predefined budget that they cannot exceed – it can be mentioned in the overall discussion of the solution

6.1.2. Non-functional

- The platform uses *OpenVZ* as a hypervisor
- The platform uses *OpenNebula* and *AppFlow* as data-center management tools
- The platform does not confine itself to one provider, but to a *ecosystem of various cloud providers* that offers deployment capabilities which vary in terms of quality of service, cost, etc.
- All communication between the user and the platform and among platform components should be encrypted

7. Evaluation

This chapter contains results and discussions on the evaluation of tests run on the proposed solution.

7.1. Introduction

We carried out a number of tests which aim to prove the solution be better than currently available, especially in terms of:

1. deployment cost
2. cost of providing given Quality-of-Service
3. deployment time

Hardware and virtual machines configuration

Many of the carried out tests were run on the same hardware and/or used the same virtual machines so for clarity of presentation their parameters are presented not in a description of every test, but here in table 7.1.

7.2. Cost of service deployment

Description

The aim of this test case is to test the primary use case of the proposed solution – deployment of a service with the emphasis of client’s **cost**. It should show that the platform chooses the best mapping between the stacks and cloud providers so that the client’s pays the **lowest** possible price.

Name	CPU	RAM	HD
(device) Desktop	Intel(R) Core(TM) i5-2500 CPU @ 3.30GHz	8GB	1TB
(device) Node1	AMD Athlon™64 X2 Dual Core, 2000MHz	3GB	160GB
(device) Node3	AMD Sempron(tm) Processor 3000+	2.5GB	160GB
(device) Node4	AMD Duron(tm) Processor 1GHz	2.5GB	60GB
(virtual machine) Frontend{1,2}	1 CPU	512MB	10GB

Table 7.1: Configuration of hardware/virtual machines used during tests

	CP-1	CP-2	CP-3
java	150	120	180
ruby	220	290	250
postgres	320	240	290
python	200	260	180
amqp	330	390	285

Table 7.2: Price for a stack in the given cloud provider

Preconditions

Service specification (A.2) forms an input to the application. Its elements are different software stacks that are parts of the whole service. Each cloud provider has its own price for a given software stack which is shown in table 7.2. Diagram 7.1 illustrates the simplified environment setup.

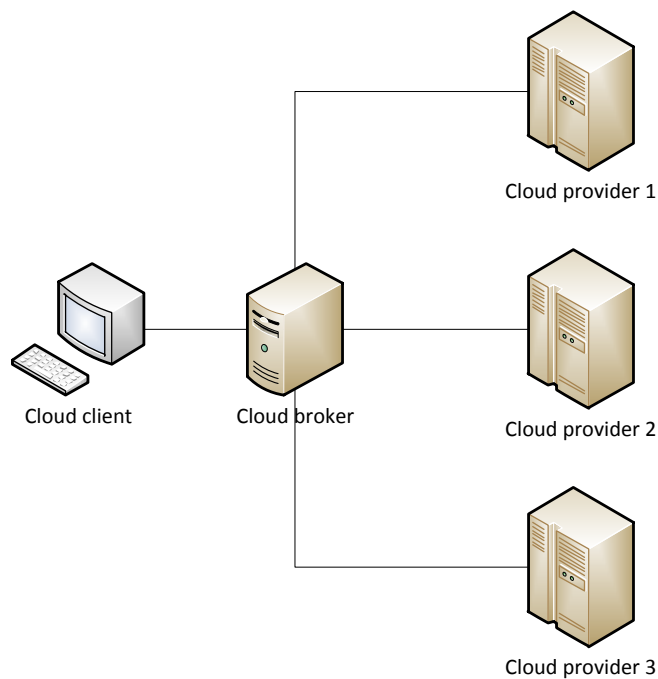


Figure 7.1: Deployment cost: environment configuration

Results

Table 7.3 shows obtained mapping between stacks and cloud providers. Taking into account this result, figure 7.2 shows comparison of cost the client would have to pay with and without such a mapping.

Conclusion

Obtained results clearly show that our proof-of-concept product met expectations of this test case as the client was offered the cheapest deployment scheme among various cloud providers for a given software stack. The map-

	CP-1	CP-2	CP-3
java		x	
ruby	x		
postgres		x	
python			x
amqp			x

Table 7.3: Chosen cloud providers for the given stack

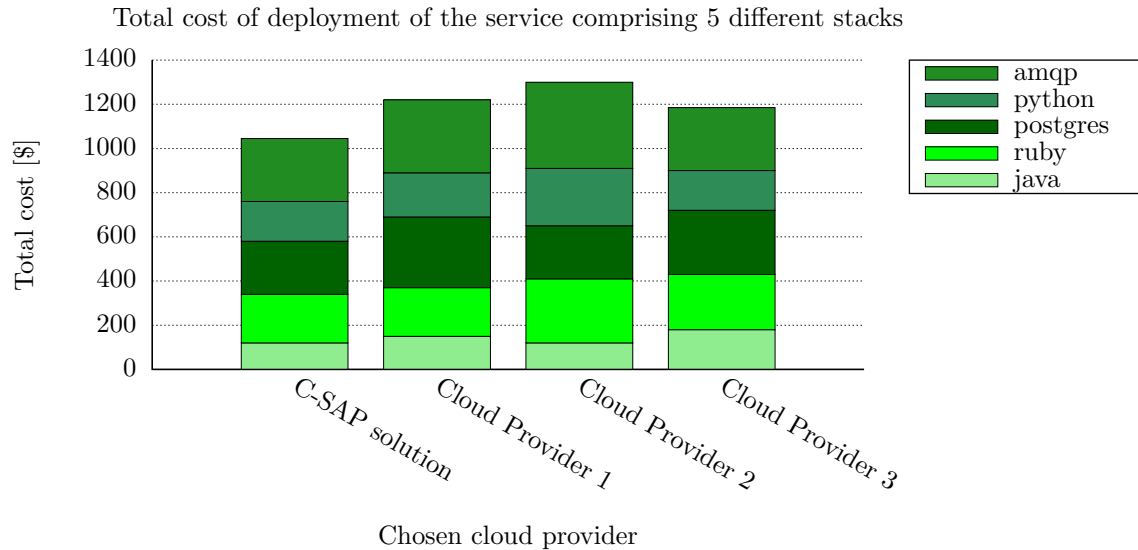


Figure 7.2: Comparison of the deployment cost when the service is deployed only on a selected cloud provider or a combination of cloud providers selected by Cloud-SAP

ping mechanism is simple yet considerably powerful – in this simple scenario the savings were significant since they constituted nearly 20 percent of the price offered by *Cloud Provider 2*. This shows that implementing similar solutions in the real world could be of great benefit to cloud consumers.

7.3. Auto-scaling – single-provider based

Description

Motivation The aim of this test-case is to show that introducing a multi-layered auto-scaling platform is of a great benefit in terms of cost and resource usage for cloud consumers and cloud providers respectively. What is more, this test should prove that once there are many levels on which scaling operations can be performed, there are significant reduction in costs paid by consumers. We want to prove it by comparing our proof-of-concept product to *Carina* [11]. *Carina* embraces a whole range of various scaling policies for users' environments (environment, in this context, means an application with all its software dependencies that are to be deployed on the cloud), but there is only one way in which they are executed – by managing the number of virtual machines (i.e. horizontal scaling). Quite on the contrary, *Cloud-SAP* has mechanisms that allow to scale the environment vertically in the first place and if that turns out to be not sufficient, horizontally. This test shows the influence of lack/presence of this feature on price and resource consumption.

Scenario The test scenario involves a) deployment of a sample environment on the cloud, b) substitute the real module responsible for collecting CPU usage data for a mock one, c) monitoring the scaling actions performed by each solution, d) evaluation of the cost the client has to pay for the service. Deployment of a service is done in a product-specific manner (up to the point where the vm deployment request is passed to OpenNebula). Mocking the CPU usage was possible by replacing the part in *InformationManager* responsible the collecting CPU data in a host for a request to a web service which generated a time-based mocked-values. Choosing the appropriate function is another issue and is discussed in the next paragraph.

CPU usage function

We wanted to ensure that both solutions, *Carina* and *Cloud-SAP*, collected the same data regarding the CPU usage for the given time. What is more, the curve should resemble CPU usage in the real business scenarios as much as possible. Thus, we took into account the following factors:

- every peak in CPU usage should be followed by a gradual descent which would mimic the real auto-scaling actions executed by each solution,
- (boundary conditions) CPU usage should be between 0 and 100 and, additionally, should exceed the previously set scaling threshold value of each solutions so that it would trigger the scaling policies evaluation
- once the scaling operations have completed, CPU usage should be constant at a rate which would not introduce any changes in the environment settings (such as the number of virtual machines and parameters of any virtual machine)

This resulted in the function whose formula is in (7.1) and plot in figure 7.3.

$$\begin{aligned}
 f(x) &= -\frac{31}{264}x^3 + \frac{3}{8}x^2 + \frac{1751}{132}x \\
 CPU_usage(t) &= \begin{cases} f(t) & 0 \leq t < 11.583 \\ f(t-10) & 11.583 \leq t < 21.441 \\ 25 & 21.441 \leq t \end{cases} \quad (7.1)
 \end{aligned}$$

Because *Carina* is capable of only horizontal scaling and *Cloud-SAP* of both horizontal and vertical, one can question whether the descent in CPU usage can be actually the same in both cases. To answer this question we want to look into the description of the virtual machines comprising the deployed environment. It states that each VM can use up to 30% power of the CPU of a host. Thus, we can be fairly sure that adding another virtual machine

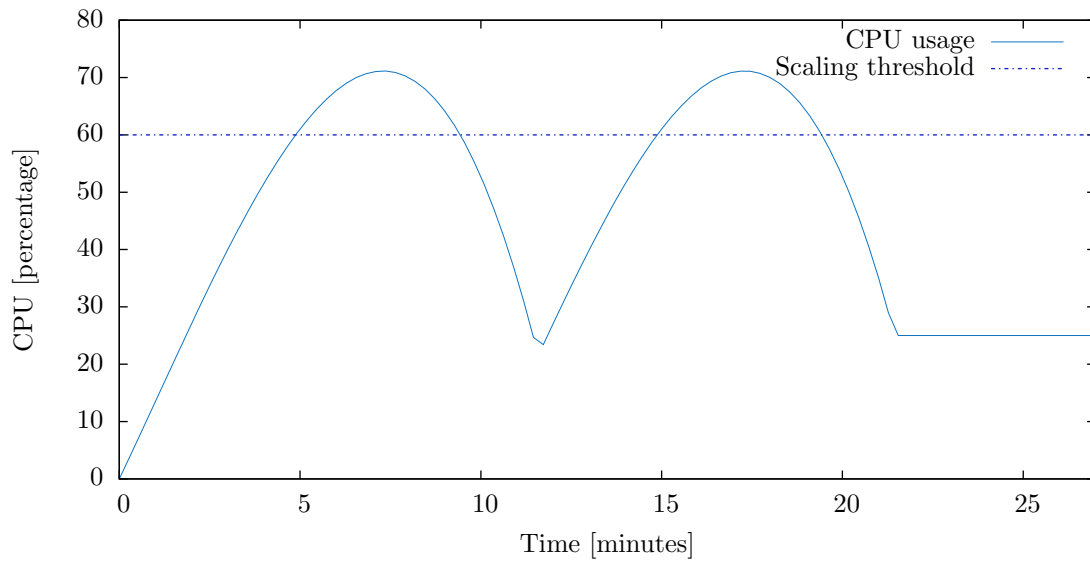


Figure 7.3: Auto-scaling - single-provider based: CPU usage function

that uses 0.3 CPU capacity of its host is equivalent with changing the CPU usage settings of already deployed VM from 30 to 60%.

Preconditions

Environment specification

Auto-scaling policy specifications As it is shown in figure 7.3, we set in both products the threshold values of CPU usage that trigger performing auto-scaling actions to 20 and 60 percent. To apply this setting, the auto-scaling part of service specification looks as follows: in Carina the user has to add the parameters of an auto-scaling policy in a hash that describes the environment under key `:elasticity_policy`. It is possible to specify the minimal and the maximal number of virtual machines that forms the environment and, what is most important, expressions which evaluation results in scaling the application. The part responsible for this is shown in listing 7.1.

Listing 7.1: Carina service specification used for testing auto-scaling with 1 cloud provider

```
ENVIRONMENT = {
  'testenv' => {
    ...
    :elasticity_policy => {
      :mode => 'auto',
      :min => 2,
      :max => 4,
      :priority => 10,
      :period => 2,
      :scaleup_expr => 'avgcpu > 60',
      :scaledown_expr => 'avgcpu < 20'
    }
  }
}
```

In Cloud-SAP it is a matter of setting appropriate arguments to a specific policy. In this case the policy is *threshold_model* and values are 20 and 60 for lower and upper bound respectively. The auto-scaling part from service specification is shown in listing 7.2.

Listing 7.2: Cloud-SAP service specification used for testing auto-scaling with 1 cloud provider

```
{
  ...
  "policies":[
    {
      "name":"threshold_model",
      "arguments": {
        "min":"20",
        "max":"60"
      }
    }
  ]
  ...
}
```

OpenNebula/Carina/Cloud-SAP configuration Since Carina and Cloud-SAP used two different instances of OpenNebula installed on separate virtual machines, it is essential the configuration be the same on each of them. Listing 7.3 shows the most important OpenNebula excerpt from settings file used in this test case – polling interval specification, which was set to 30 seconds. To ensure that both products can actually use up-to-date data, they should evaluate their policies every 30 seconds or longer. In Carina the user cannot specify this value, because it is hard-coded in source code to 60 seconds. In Cloud-SAP, however, this value was set in a configuration file and was equal to 60 seconds as well as in Carina.

Listing 7.3: OpenNebula configuration excerpt – information manager

```
HOST_MONITORING_INTERVAL = 300
HOST_PER_INTERVAL        = 15
VM_POLLING_INTERVAL       = 30
VM_PER_INTERVAL           = 10
```

. Vertical scaling mechanism in *Cloud-SAP* is implemented in a way that causes exponential growth of CPU resources consumption. When there is a need to perform a scaling action on CPU usage, current value is taken and increased by 30%.

Deployment diagrams Deployment diagrams show the placement of each component in both installations. They are shown in figures 7.4 and 7.5.

Hardware/VM configuration

All virtual machines were deployed onto *Node1*, which uses *OpenVZ* as a virtualization technology and whose hardware configuration can be found in table 7.1. *OpenNebula* was installed on a *Frontend1* virtual machine and its configuration is in the same table. *Carina* required the usage of 2 virtual machines which were deployed on *Desktop* with KVM as a hypervisor and whose configuration is in table 7.4. Diagram 7.6 illustrates the physical setup of the test-case.

Expectations

As auto-scaling behaviour of both products is determined by CPU usage of the deployed service, we can predict the outcome of the test by thorough analysis of the plot in figure 7.3.

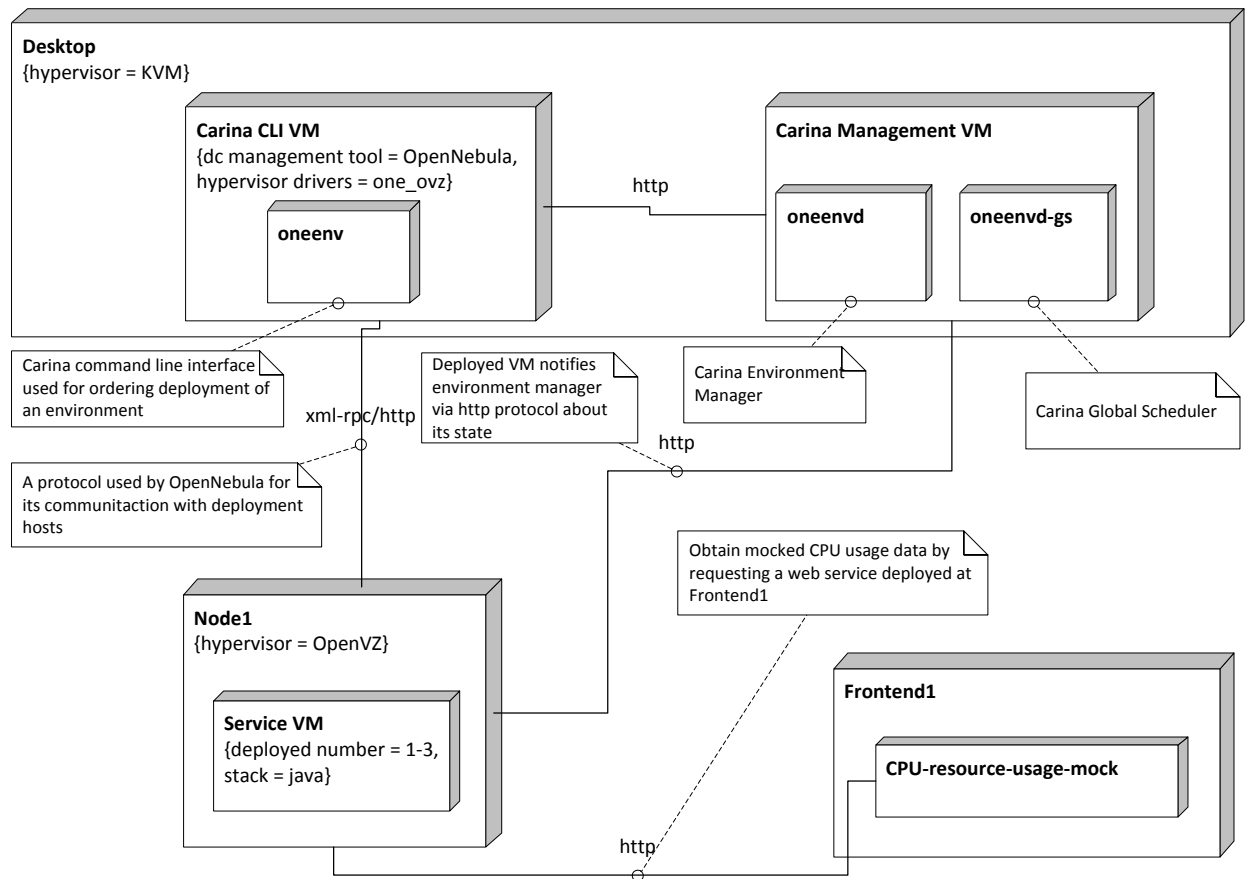


Figure 7.4: Auto-scaling - single-provider based: deployment diagram of Carina

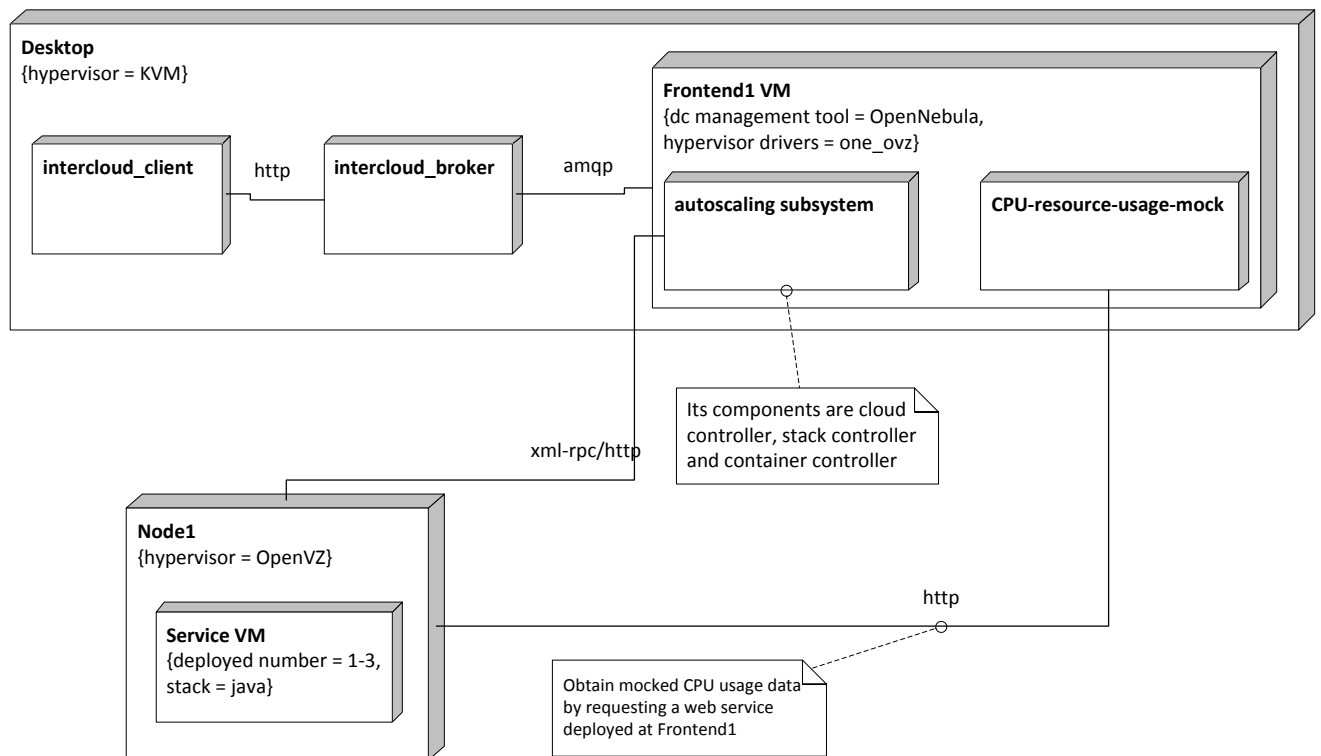


Figure 7.5: Auto-scaling - single-provider based: deployment diagram of Cloud-SAP

Name	CPU	RAM	HD
(virtual machine) carina-frontend	1 CPU	1GB	6GB
(virtual machine) carina-management-vm	1 CPU	1GB	9GB

Table 7.4: Configuration of hardware/virtual machines used during tests

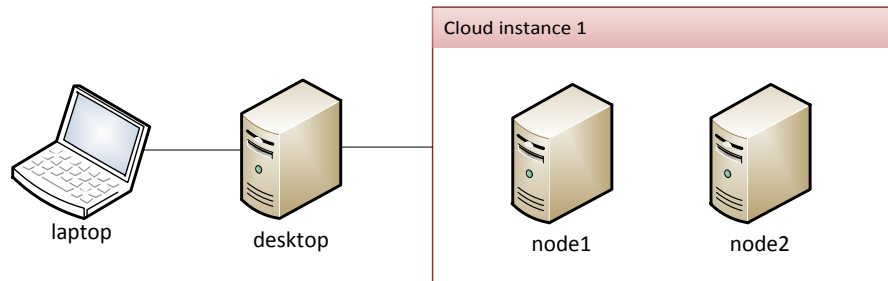


Figure 7.6: Auto-scaling - single-provider based: environment configuration

The first conclusion is that all auto-scaling actions are performed after exceeding the upper threshold value of CPU usage, this is 60 percent in our case what happens about 5 minutes after the deployment of an environment is completed. CPU consumption remains high for about 2 minutes what enables both products to take appropriate, auto-scaling steps. We expect *Carina* to deploy another virtual machine and *Cloud-SAP* to scale vertically. Then it gradually goes down, which simulates that all taken actions have successfully completed. Then, for another several minutes, the cycle recurs, but this time we expect the proposed solution to scale horizontally. *Carina* is expected to behave as in the previous cycle. Once the cycle has completed, after roughly 21 minutes, the value of CPU consumption remains constant at 25 percent and the test is completed.

Results

In the first place it is worth discussing the results of *Carina* and *Cloud-SAP* separately and then compare resource consumption of each solution and its influence on cost paid by the end-user.

Cloud-SAP Our proof-of-concept solution behaved exactly as expected. To make the discussion more clear, there is an excerpt from a log file of this test run in listing 7.4.

Listing 7.4: Cloud-SAP logs excerpt regarding auto-scaling actions

```

...
[2013-12-31T15:25:45.05] DEBUG : Notifying process 2736 about the deployed service
...
[2013-12-31T15:31:45.98] DEBUG : CONTAINER Concluded that currently {"id":1,"correlation_id":
    "875","ip":"192.168.0.100","type":"master","probed":"1388500305","requirements":{"cpu":
    "0.3","memory":512},"stack_id":1} is insufficient (by key: CPU)
[2013-12-31T15:31:45.98] DEBUG : CONTAINER Attempt to scale CPU up for a container: {"id":
    "1","correlation_id":875,"ip":"192.168.0.100","type":"master","probed":"1388500305","
    requirements":{"cpu":0.3,"memory":512},"stack_id":1}
D, [2013-12-31T15:31:46.78] DEBUG : Prepared payload for CPU increase: {"cpulimit"=>39} for
    {"id":1,"correlation_id":875,"ip":"192.168.0.100","type":"master","probed":
    "1388500305","requirements":{"cpu":0.39,"memory":512},"stack_id":1} at node1
...

```

```
[2013-12-31T15:32:46.15] DEBUG : CONTAINER Attempt to scale CPU up for a container: {"id":1,"correlation_id":875,"ip":"192.168.0.100","type":"master","probed":"1388500365","requirements":{"cpu":0.39,"memory":512},"stack_id":1}
[2013-12-31T15:32:46.94] DEBUG : Prepared payload for CPU increase: {"cpulimit"=>50} for {"id":1,"correlation_id":875,"ip":"192.168.0.100","type":"master","probed":"1388500365","requirements":{"cpu":0.507,"memory":512},"stack_id":1} at node1
...
[2013-12-31T15:33:46.32] DEBUG : CONTAINER Concluded that currently {"id":1,"correlation_id":875,"ip":"192.168.0.100","type":"master","probed":"1388500425","requirements":{"cpu":0.507,"memory":512},"stack_id":1} is insufficient (by key: CPU)
D, [2013-12-31T15:33:47.12] DEBUG : Prepared payload for CPU increase: {"cpulimit"=>65} for {"id":1,"correlation_id":875,"ip":"192.168.0.100","type":"master","probed":"1388500425","requirements":{"cpu":0.6591,"memory":512},"stack_id":1} at node1
...
[2013-12-31T15:43:50.27] INFO : STACK Delegating execution to a cloud-controller
[2013-12-31T15:43:50.28] INFO : Received request of insufficient_slaves to be performed on a stack #<AutoScaling::Stack @id=1 @correlation_id =628 @type=:java @state=:deployed @data=nil @service_name="Auto-scaling test">
...
```

Carina As *Carina* is capable only of horizontal scaling, in discussion of its resource usage we can confine ourselves to counting the number of deployed virtual machines in the given time intervals. In the log files we can trace all actions that were triggered during the test case. As listing 7.5 shows, they completely met the expectations expressed in the previous section – there were two “SCALEUP” jobs which resulted in increase of the total number of virtual machines comprising the environment.

Listing 7.5: Carina environment manager logs with taken actions (*jobs*)

ID	ENVID	CONFIG_NAME	TYPE	SUBMIT_TIME	STATUS
--	-----	-----	----	-----	-----
206	12	testenv	CREATE	Tue Dec 31 14:38:23 +0000 2013	DONE
207	12	testenv	SCALEUP	Tue Dec 31 14:48:07 +0000 2013	DONE
208	12	testenv	SCALEUP	Tue Dec 31 14:58:07 +0000 2013	DONE

. Since the base configuration assumed that the environment consists of 2 virtual machines, we can say that about 10 minutes after deployment of an environment, the number of VMs increased to 3, and after another 10 minutes, to 4.

Deployment finished at 15:25:45 and at this time we started the mock-cpu-usage web service. When the platform concluded that there are insufficient resources, it tried for three times changing parameters (virtual CPU used by the vm) of the virtual machines forming the service. Once scaling vertically was not possible such an information was passed to cloud controller which performed scaling across different cloud providers.

Cost

In our scenario cost is roughly equivalent with the CPU usage, so the plot of a cost virtually presents the cpu usage. Figure 7.7 shows cpu usage by environment when deployed and configured by two competing products. If we were to roughly estimate the cost in both solutions it would be equal in Carina:

$$7 \cdot 0.3 + 10 \cdot 0.6 + 8 \cdot 0.9 = 15.3 \quad [\text{currency unit}] \quad (7.2)$$

and in Cloud-SAP:

$$7.0166 \cdot 0.3 + 2 \cdot 0.39 + 2.01667 \cdot 0.507 + 7.05 \cdot 0.659 + 6.91667 \cdot 0.959 \approx 15.186483 \quad (7.3)$$

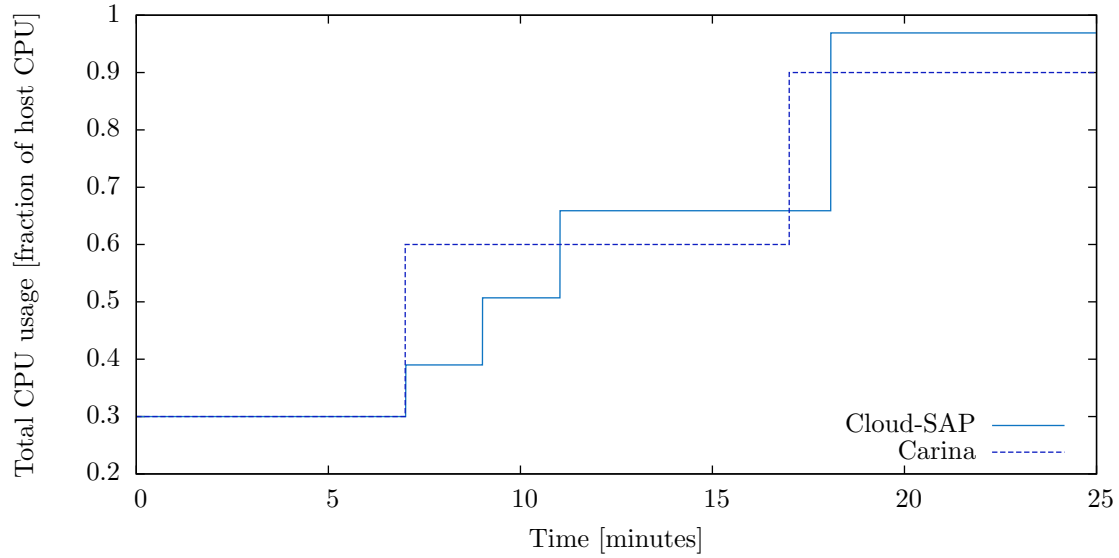


Figure 7.7: Auto-scaling - single-provider based: comparison of cost/CPU usage

Conclusion

The obtained results shows that the proposed solution has enormous potential in terms of better utilization of available resources of the cloud and reducing cost paid by cloud consumers.

Judging by the plot of cost (figure 7.7), we can be fairly sure that it is possible to obtain better results by proper tuning of the scaling vertical mechanism. In particular, one can consider changing the default CPU growth from 1.3 to other or apply a different strategy, for example a growth by a constant rate. What is more, in order to get the best results, the policy scaling interval must be chosen with a great care with the consideration of specific parameters of a given environment.

7.4. Auto-scaling – multiple-provider based

Description

Motivation

This test case aims to prove that scaling across multiple cloud providers is vitally important while ensuring appropriate Quality-of-Service, especially in cases of increased number of service requests. Such scaling scenario, known also as cloud-bursting, leverage benefits arising from offloading application load to an external provider. In order to verify our concept we compare number of transactions per second guaranteed by *Cloud-SAP* and *Carina*. While our proof of concept solution features multi-layer scaling and is cloud federation aware, *Carina* adopts only horizontal scaling. It is expected that test case prove benefits emerging from enriching application platform provider with an cloud federation awareness.

Scenario

Testing requires following steps to be done:

1. deploy exemplary service
2. observe system behaviour under load, simulated by a resource usage mock. It is vital for a resource usage to exceed a single provider capabilities
3. assess system characteristics, including number of handled transactions per second

Load simulation

Similarly to a previous test case, we leveraged a resource mock that allows us to precisely control monitoring information returned to an system. Besides, in order to simplify test case, we were solely focused on a CPU usage. CPU usage function has to exceed upper threshold limit at some point and stay at that level, triggering successive auto-scaling events. However, at some point single cloud provider resources will be surpassed. Equation (7.4) denotes a resource usage function that is expected to fulfil above-mentioned scenario and is illustrated in figure 7.8.

$$f(x) = \frac{29}{1500}x^3 - \frac{21}{20}x^2 + \frac{533}{30}x \quad (7.4)$$

Preconditions

Environment specification

Auto-scaling policy specifications Policies used in this test are based on a threshold model with the following properties:

- value below 20 triggers scaling down event
- value grater than 60 triggers scaling up event

However, considering the fact that we are entirely focused on scaling up, only the upper limit is relevant in our case. Listings 7.6 and 7.7 presents scaling policies for *Carina* and *Cloud-SAP* respectively.

Listing 7.6: *Carina* service specification used for testing auto-scaling with 2 cloud providers

```
ENVIRONMENT = {
  'testenv' => {
    ...
```

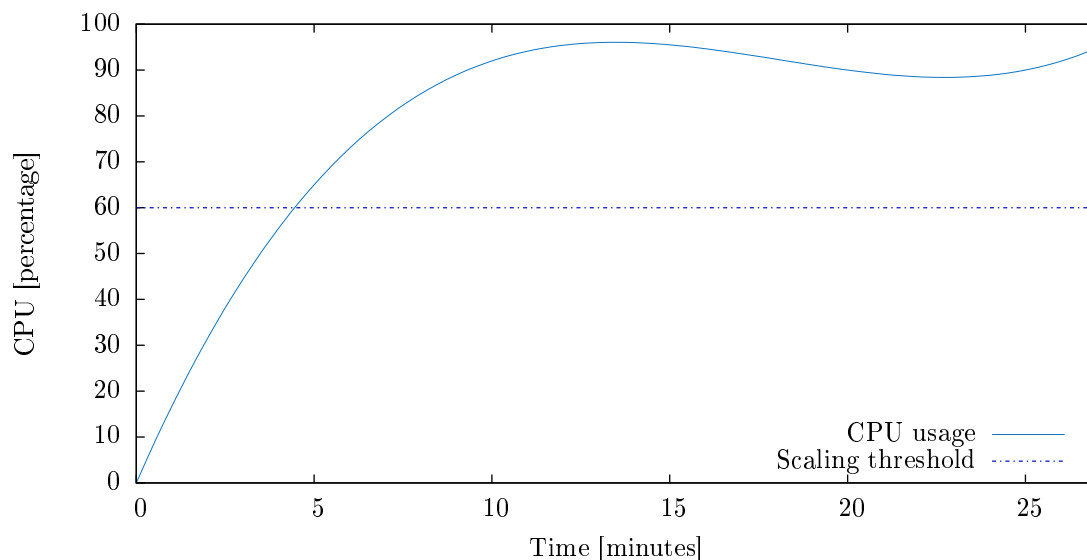


Figure 7.8: Auto-scaling - multiple-provider based: CPU usage function

```

:elasticity_policy => {
  :mode => 'auto',
  :min => 2,
  :max => 10,
  :priority => 10,
  :period => 2,
  :scaleup_expr => 'avgcpu > 60',
  :scaledown_expr => 'avgcpu < 20'
}
}
}

```

Listing 7.7: Cloud-SAP service specification used for testing auto-scaling with 2 cloud providers

```

{
  ...
  "policies":[
    {
      "name":"threshold_model",
      "arguments": {
        "min":"20",
        "max":"60"
      }
    }
  ]
  ...
}

```

OpenNebula/Carina/Cloud-SAP configuration Components common to *Carina* and *Cloud-SAP*, namely *OpenNebula* instances and computing nodes, were configured in the same fashion. Listing 7.8 depicts key configuration elements of *OpenNebula* monitoring mechanism.

Listing 7.8: OpenNebula configuration excerpt – virtual machine and information manager

```

HOST_MONITORING_INTERVAL = 300
HOST_PER_INTERVAL        = 15
VM_POLLING_INTERVAL      = 30
VM_PER_INTERVAL          = 10

```

Deployment diagrams Components that took part during testing are show in figures 7.10 and 7.9.

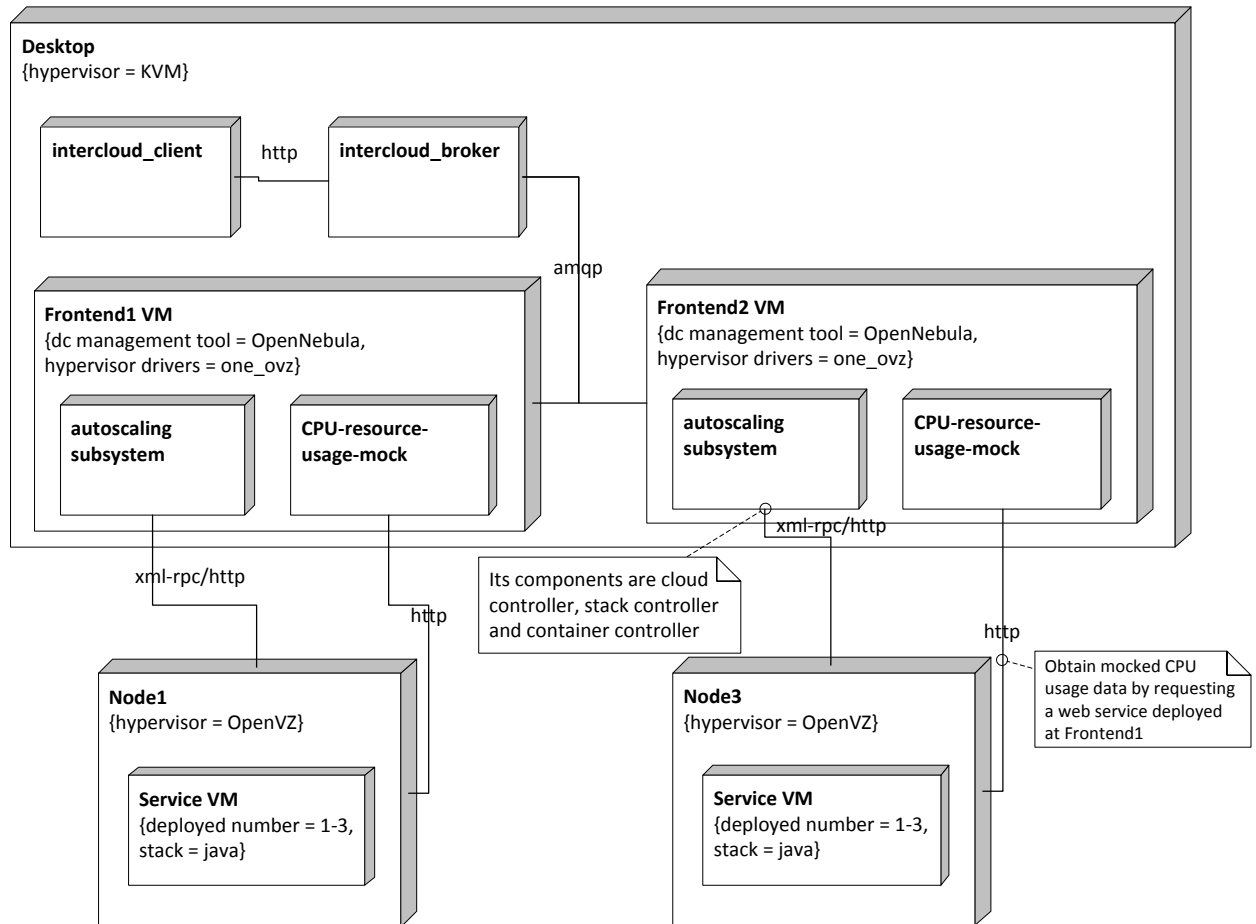


Figure 7.9: Auto-scaling - multiple-provider based: deployment diagram of Carina

Hardware/VM configuration

Figure 7.11 depicts physical configuration of the environment. In short, setup was as follows: all *Cloud-SAP* components, apart from cloud client deployed on *Laptop*, were provisioned on *Desktop*. *Cloud Provider 1 (CP-1)* is *OpenNebula* instance known as *Frontend1* which uses *Node1* as a computing node and is deployed on *Desktop*, while *CP2* uses *Frontend2* and *Node2*. Specification of nodes is listed in table: 7.1. Deployment cost of a java stack was 50 and 60 using *CP-1* and *CP-2*, respectively.

Expectations

Taking CPU usage function into account, we can predict outcome of the test. Scaling policy is violated in 5th minute, hence, it is expected that first scaling event occur after that point. Moreover, since that moment, CPU usage remains beyond threshold implying successive scaling actions.

Due to the fact that *Cloud-SAP* favours fine-grained scaling actions such as vertical scaling, it is expected that these actions occur sooner than horizontal scaling. However, at some point *CP-1* capacity won't be sufficient,

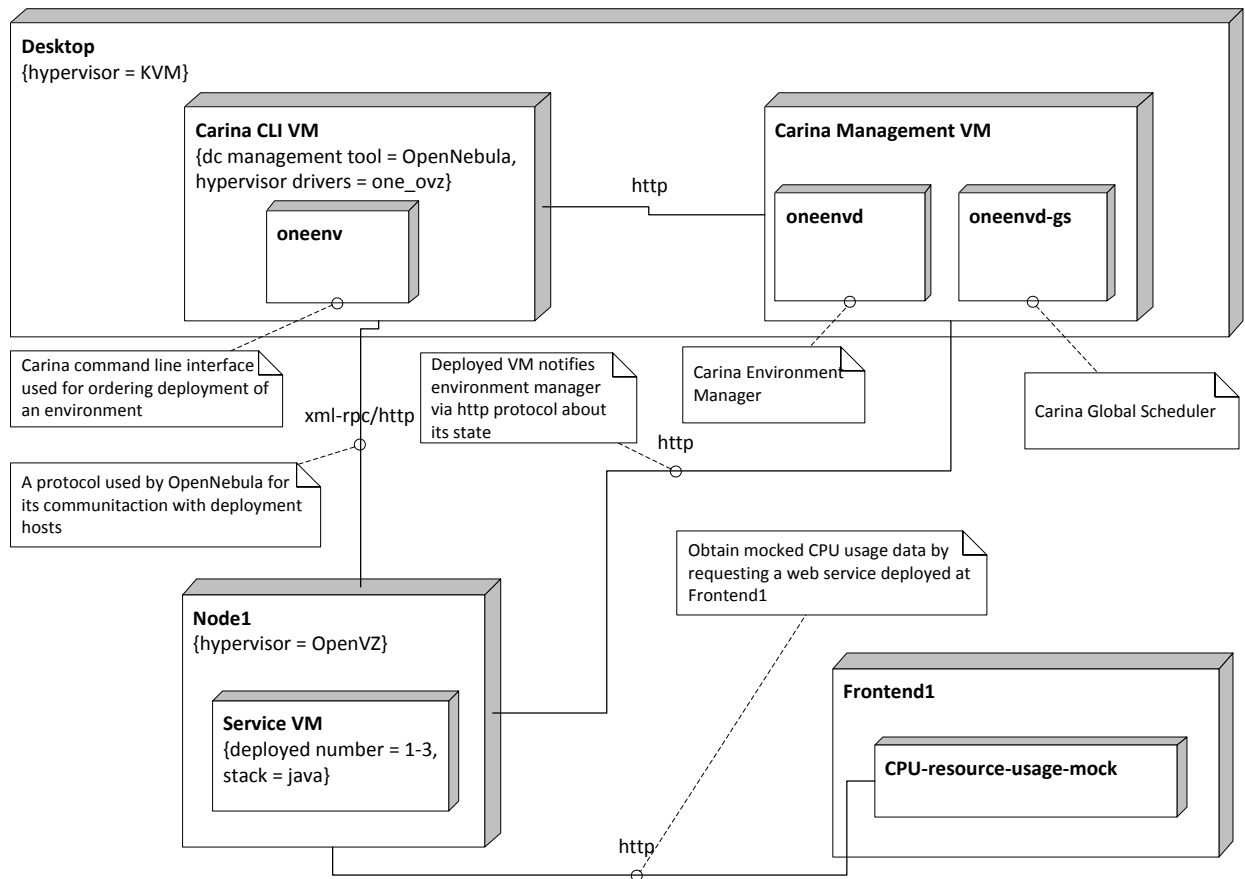


Figure 7.10: Auto-scaling - multiple-provider based: deployment diagram of Carina

hence, another stack instance (one master instance, two slaves) will be deployed on *CP-2*. Subsequent scaling actions will take place solely on a *CP-2* 2 and will involve further vertical scaling.

Carina, on the other hand, supports solely horizontal scaling using single cloud provider. On top of that, it is assumed that succeeding slaves instances will be added to a service up to the point where Cloud Provider won't have enough resources to proceed with further scaling requests.

Results

Cloud-SAP In total, there were 13 vertical (increasing CPU limit) 2 horizontal (adding a container) scaling events. First action took place minute after first violation of scaling policy rule. Since then, slaves' CPU were successively increasing to the point where further scaling wasn't possible - 13th minute of the test. Therefore, *Cloud-SAP* deployed two new slaves on *CP-2*. Listing 7.9 presents logs covering service lifecycle.

Listing 7.9: Cloud-SAP logs excerpt regarding auto-scaling actions on multiple providers

```
D, [2014-01-03T15:43:48.782879 #2740] DEBUG -- : CONTAINER Analyzing data {container=>#<
AutoScaling::Container @id=1 @correlation_id=875 @ip=#<IPAddr: IPv4
:192.168.0.100/255.255.255.255> @type=:master @probed="1388501025" @requirements={"cpu
"=>1.8824555100000004, "memory"=>512} @stack_id=1>, :metrics=>{"CPU"=>["71", "70"]}}
D, [2014-01-03T15:43:48.783877 #2740] DEBUG -- : CONTAINER Concluded that currently {"id
":1,"correlation_id":875,"ip":"192.168.0.100","type":"master","probed":"1388501025","
requirements":{"cpu":1.8824555100000004,"memory":512},"stack_id":1} is insufficient (by
key: CPU)
...
```

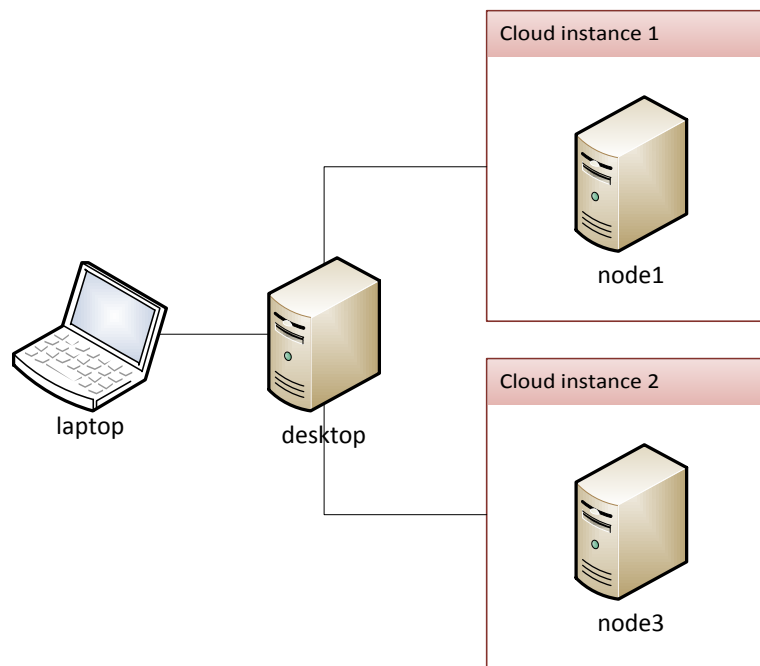


Figure 7.11: Auto-scaling - multiple-provider based: environment configuration

```

D, [2014-01-03T15:43:48.784777 #2740] DEBUG -- : CONTAINER Attempt to scale CPU up for a
  container: {"id":1,"correlation_id":875,"ip":"192.168.0.100","type":"master","probed
    ":"1388501025","requirements":{"cpu":1.8824555100000004,"memory":512},"stack_id":1}
I, [2014-01-03T15:43:49.512883 #2740] INFO -- : CONTAINER Cannot reserve:
  2.447192163000001 with {:memory=>358.76171875, :cpu=>1.9480000000000002} at node1 (
  AutoScaling::InsufficientResources)
...
I, [2014-01-03T15:43:49.514541 #2740] INFO -- : STACK Got unprocessed conclusion:
  insufficient_cpu for {"id":2,"correlation_id":876,"ip":"192.168.0.101","type":"slave",
  "probed":"1388501010","requirements":{"cpu":1.8824555100000004,"memory":512},"stack_id
    ":1}
...
I, [2014-01-03T15:43:50.279589 #2740] INFO -- : STACK Cannot reserve: {:cpu
  =>1.8824555100000004, :memory=>512.0} with {:cpu=>0.3675444899999998, :memory
  =>2673.109375} (AutoScaling::InsufficientResources)
I, [2014-01-03T15:43:50.279955 #2740] INFO -- : STACK Delegating execution to a cloud-
  controller
I, [2014-01-03T15:43:50.280566 #2740] INFO -- : Received request of insufficient_slaves to
  be performed on a stack #<AutoScaling::Stack @id=1 @correlation_id=628 @type=:java
  @state=:deployed @data=nil @service_name="Deployment time test service">
  
```

Carina Similarly to a *Cloud-SAP*, first scaling event (slave addition) was noticed in 6th minute of test. Scaling jobs were continuously invoked every 2 minutes until capacity of a *CPI* has been fully exploited. Listing 7.10 summarises jobs performed by *Carina*.

Listing 7.10: Carina environment manager logs with taken actions (*jobs*)

ID	ENVID	CONFIG_NAME	TYPE	SUBMIT_TIME	STATUS
--	-----	-----	----	-----	-----

266	17	testcase2	CREATE	Fri Jan 3 10:41:15 +0000 2014	DONE
267	17	testcase2	SCALEUP	Fri Jan 3 10:47:44 +0000 2014	DONE
268	17	testcase2	SCALEUP	Fri Jan 3 10:49:31 +0000 2014	DONE
267	17	testcase2	SCALEUP	Fri Jan 3 10:51:02 +0000 2014	DONE
268	17	testcase2	SCALEUP	Fri Jan 3 10:53:07 +0000 2014	DONE

Transactions per second

Chart 7.12 illustrates how service capabilities changed over the time. Service capability is expressed as a number of transactions per second that can be handled by a service. Please note, that this measure, while giving good insight into application potential, is hard to simulate, hence, we assumed that 1 VCPU provides enough resources to successfully serve 100 TPS. Knowing that, we roughly estimate total service capacity, which at the end of the test amounts to:

- Cloud-SAP: 283 transactions per second
- Carina: 180 transactions per second

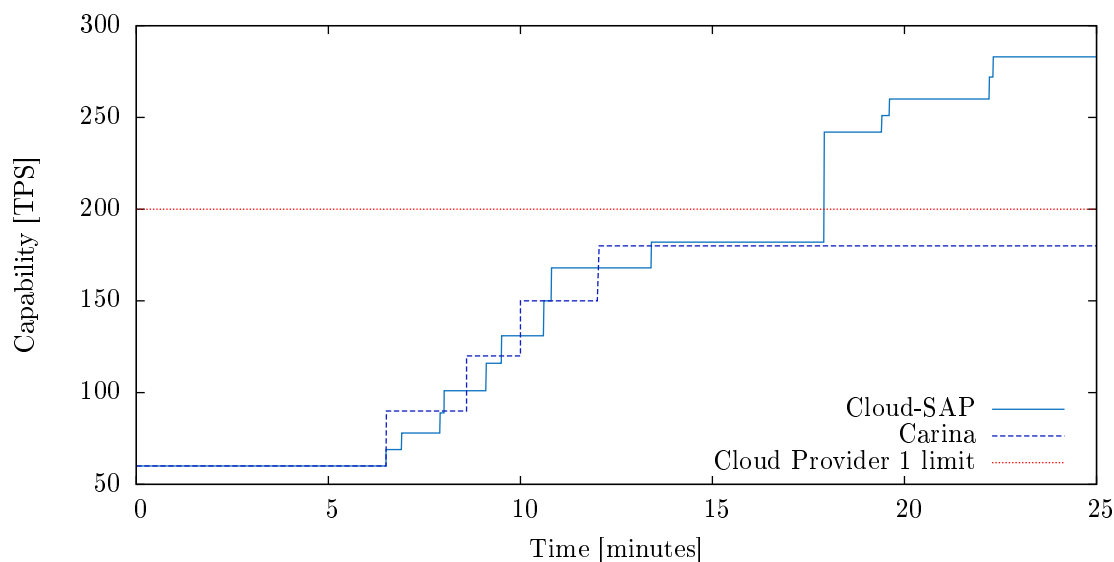


Figure 7.12: Auto-scaling - multiple-provider based: comparison of processed transaction per second

Conclusion

Comparing expectations with test output, one can see that expectations has been fully met. While *Carina* was solely focused on horizontal scaling and bound to a single cloud provider, *Cloud-SAP* first took fine-grained actions to later leverage resources of a second cloud provider.

Cloud-SAP has been proven to be a better solution, offering greater capabilities to a service. This was possible by exploiting resources of two federated cloud providers, crucial concept behind *Cloud-SAP* design.

7.5. Deployment time – solution comparison

Description

In this test we want to compare our solution to one of those available at the market which use *OpenNebula* as an underlying tool for managing resources of a data center and *OpenVZ* as a hypervisor in terms of **deployment time**, one of the most important factors of products whose main purpose is to scale applications. *Carina* [11] can be considered a perfect match of a solution for such a comparison and tests are ran against it.

This test involves the steps of i) instantiating one of the tested product, i.e. Cloud-SAP or Carina, ii) ordering the deployment of a service whose specification is shown in listing A.3, iii) measuring the time needed to set up the environment of the service.

Preconditions

It is assumed that *Cloud-SAP* and *Carina* according with *OpenVZ* as an underlying virtualization technology are correctly installed and configured. Each test case must be run in an isolation so before performing any test all virtual machines present at the deployment node are removed.

OpenNebula configuration

To ensure objectivity in tests, OpenNebula was configured in both products in the same way. One of the key factors that could influence the deployment time is the configuration of scheduler. Its parameters are shown in the listing 7.11.

Listing 7.11: OpenNebula scheduler configuration

```
SCHED_INTERVAL = 30
MAX_VM          = 300
MAX_DISPATCH    = 30
MAX_HOST        = 1
HYPERVISOR_MEM  = 0.1
```

Service description

The service comprises a simple java enterprise application, deployed in a master-slave configuration with one VM set as a load balancer and other nodes that serve as workers, which uses Tomcat as a web container.

The description of a service expressed in Carina format can be found in listing A.1.

Hardware configuration

All virtual machines were deployed on a host named *Node1*, whose configuration can be found in table 7.1. Diagram presents the physical configuration of nodes.

Environment configuration

Deployment diagram for *Carina* implementation is shown in figure 7.14 and for *Cloud-SAP* in figure 7.15.

Results

Obtained results are shown in table 7.5 and in figure 7.16. For a given number of instances we ordered deploying a service 10 times and the values shown in those figures are an average of these runs.

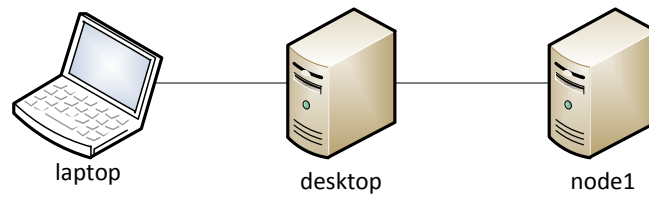
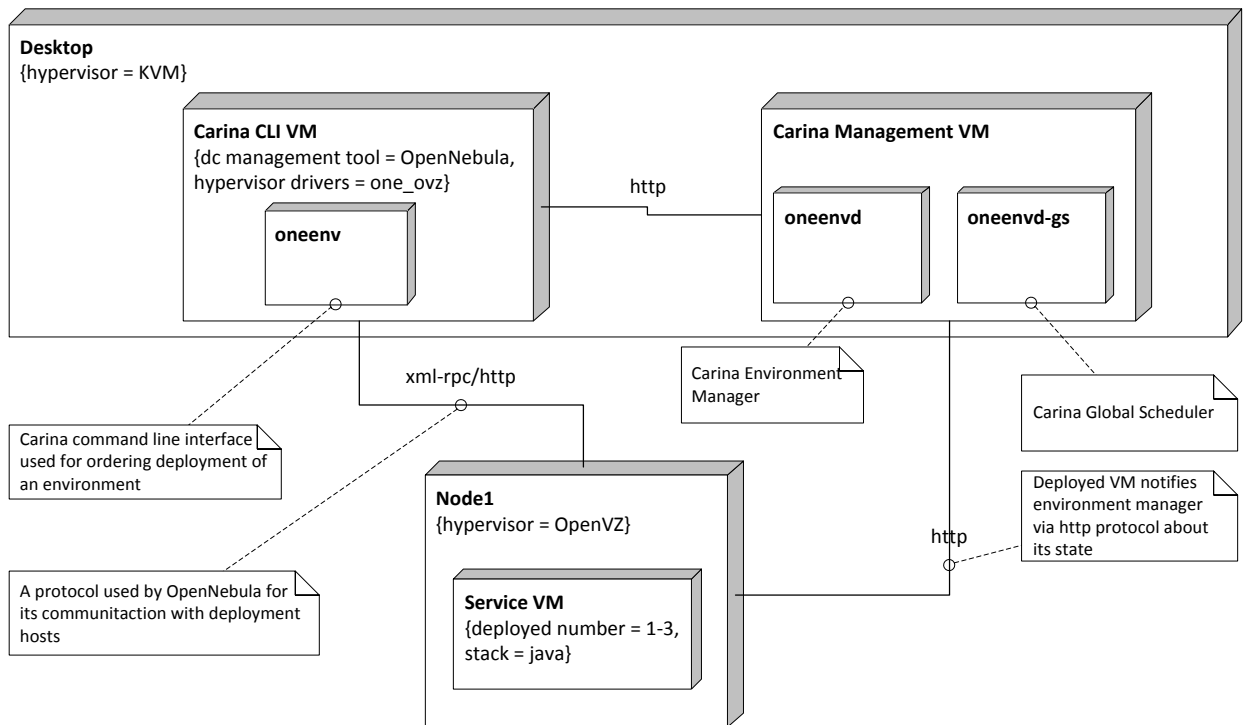


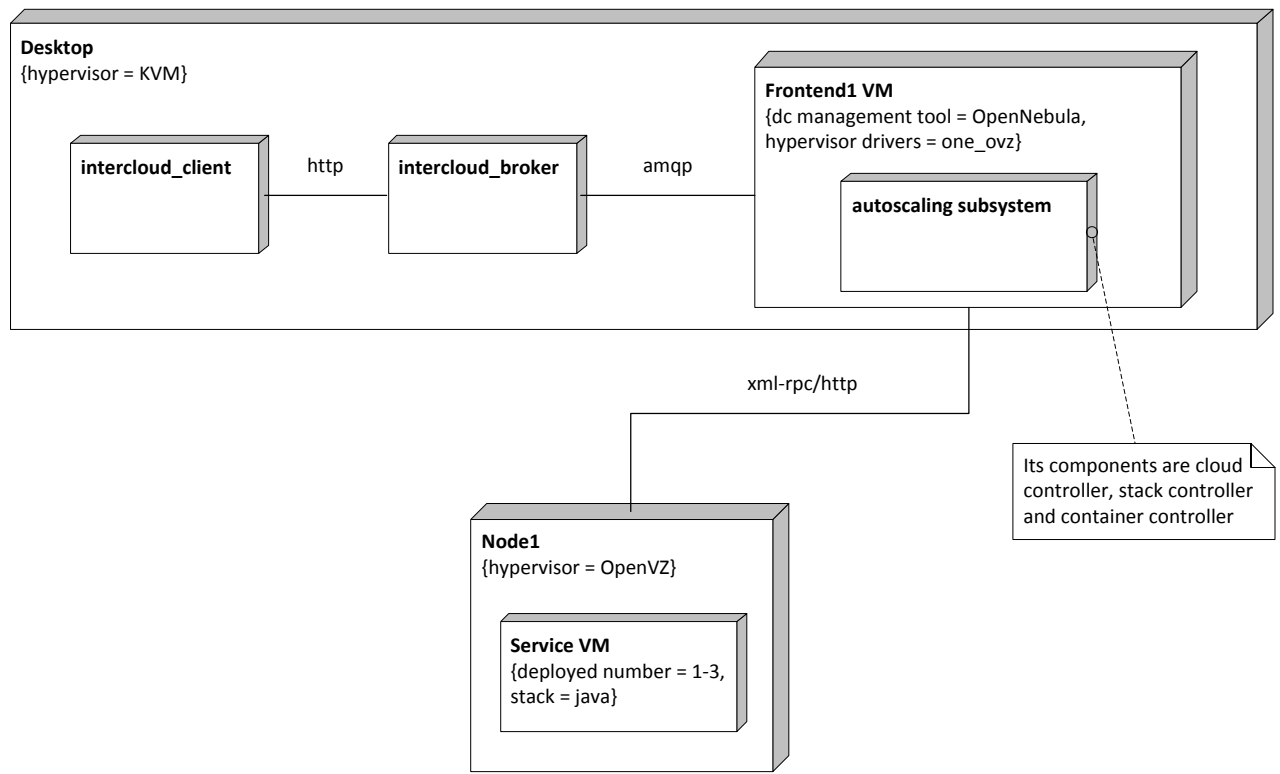
Figure 7.13: Deployment time - solution comparison: physical environment setup

Figure 7.14: Deployment diagram of *Carina*

Conclusion

As one can notice, deployment of a java stack takes using *Cloud-SAP* lasts longer than when using *Carina* for the same purpose. Moreover, the more instances are in such stack, the greater difference is. This difference in provisioning time is mainly driven by the fact that *Cloud-SAP* uses a great deal of components in comparison to *Carina*, what is a direct consequence of chosen architecture. Specifically, the request from a client is passed to a broker, then the best provider is chosen and finally the deployment using the select provider and its *Applfow* server takes place. Contrary, *Carina* directly operates on *OpenNebula* giving much simpler flow.

Noticeably, deployment time never was a crucial issue for *Cloud-SAP*, hence, this issue is of a little importance. What is actually important for *Cloud-SAP* is deployment cost and service performance. In other words, slightly greater deployment time is a price that is paid to be sure that resources are deployed optimally and with Quality-of-Service guarantee.

Figure 7.15: Deployment diagram of *Cloud-SAP*

Instance no	Solution	
	Cloud-SAP	Carina
2	198.0	158.7
3	261.72	208.1
4	294.38	230.6

Table 7.5: Average deployment time for the service with the various number of VMs used for the whole environment

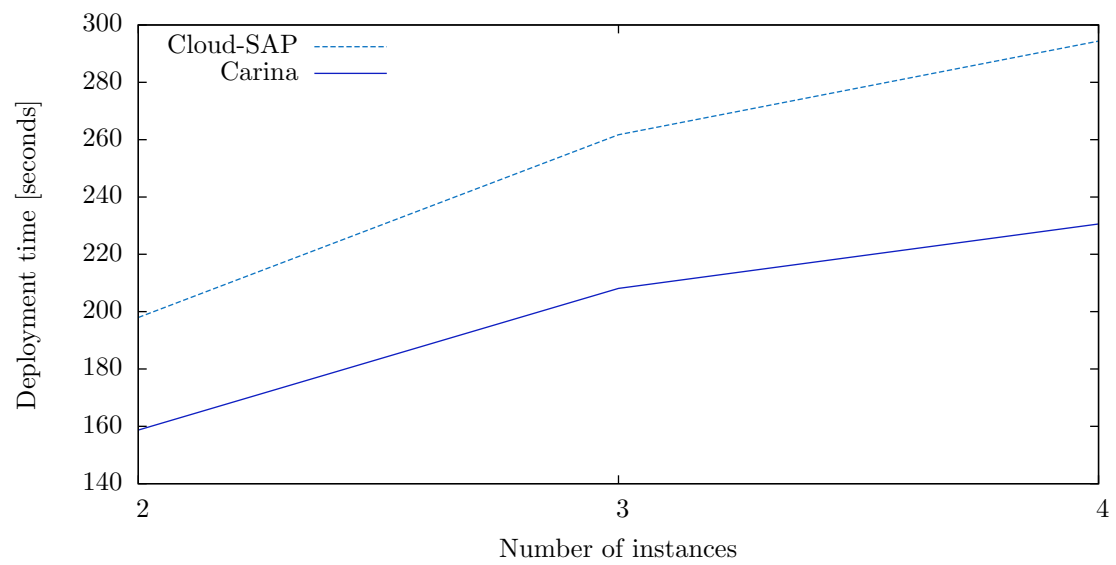


Figure 7.16: Average deployment time for two competing products when the variable is the number of instances of VMs

8. Summary

A. Code listings

A.1. Service specifications

Listing A.1: Carina Environment Specification which was used during tests of deployment time

```
ENDPOINT = {
  'mm01' => {
    :proxy    => 'http://192.168.0.35:2633/RPC2',
    :oneauth => "svc:xxxxx"
  }
}

TEMPLATE = {
  'tomcat' => {
    :file      => "~/vm/tomcat.vm",
    :cpu       => "0.3",
    :memory    => 512,
    :network_id => { 'mm01' => 7 },
    :image_id  => { 'mm01' => 10 }
  },
  'haproxy' => {
    :file      => "~/vm/haproxy.vm",
    :cpu       => "0.3",
    :memory    => 512,
    :network_id => { 'mm01' => 7 },
    :image_id  => { 'mm01' => 11 }
  }
}

ENVIRONMENT = {
  'testenv' => {
    :type                => "compute",
    :endpoint            => "mm01",
    :description         => "Example environment",
    :master_template     => "haproxy",
    :master_context_script => "sample-master-context-script.sh",
    :master_setup_time   => 30,
    :master_context_var  => "BALANCE_PORT=8080",
    :slave_template      => "tomcat",
    :slave_context_script => "sample-slave-context-script.sh",
    :slave_context_var   => "APP_PACKAGE=gwt-petstore.war",
    :placement_policy    => "pack",
    :num_slaves          => 3,
  }
}
```

```

        :slavedata           => "8080",
        :adminuser           => "root",
        :app_url              => "http://%MASTER%:8080/testapp"
    }
}

```

Listing A.2: Cloud-SAP Service Specification (without scaling policies)

```

{
  "name": "my new facebook",
  "stacks": [
    {
      "type": "java",
      "instances": 1
    },
    {
      "type": "amqp",
      "instances": 1
    },
    {
      "type": "python",
      "instances": 1
    },
    {
      "type": "ruby",
      "instances": 1
    },
    {
      "type": "postgres",
      "instances": 1
    }
  ]
}

```

Listing A.3: Cloud-SAP Service Specification used for testing deployment time

```

{
  "name": "Deployment time test service",
  "stacks": [
    {
      "type": "java",
      "instances": 2,
      "policy_set": {
        "min_vms": 0,
        "max_vms": 2,
        "policies": [
          {
            "name": "threshold_model",
            "parameters": {
              "min": "5",
              "max": "50"
            }
          }
        ]
      }
    }
  ]
}

```



```

    ]
}

```

A.2. Scaling policies

Listings A.4 illustrates XML-based policy that is used by Auto Scaling of the Amazon Web Services EC2.

Listing A.4: Scaling policy - AWS EC2

```

<DescribeAutoScalingGroupsResponse xmlns="http://autoscaling.amazonaws.com/doc
/2011-01-01/">
  <DescribeAutoScalingGroupsResult>
    <AutoScalingGroups>
      <member>
        <Tags/>
        <SuspendedProcesses/>
        <AutoScalingGroupName>my-test-asg</AutoScalingGroupName>
        <HealthCheckType>EC2</HealthCheckType>
        <CreatedTime>2013-01-22T23:58:48.718Z</CreatedTime>
        <EnabledMetrics/>
        <LaunchConfigurationName>my-test-lc</LaunchConfigurationName>
        <Instances>
          <member>
            <HealthStatus>Healthy</HealthStatus>
            <AvailabilityZone>us-east-1e</AvailabilityZone>
            <InstanceId>i-98e204e8</InstanceId>
            <LaunchConfigurationName>my-test-lc</LaunchConfigurationName>
            <LifecycleState>InService</LifecycleState>
          </member>
        </Instances>
        <DesiredCapacity>1</DesiredCapacity>
        <AvailabilityZones>
          <member>us-east-1e</member>
        </AvailabilityZones>
        <LoadBalancerNames/>
        <MinSize>1</MinSize>
        <VPCZoneIdentifier/>
        <HealthCheckGracePeriod>0</HealthCheckGracePeriod>
        <DefaultCooldown>300</DefaultCooldown>
        <AutoScalingGroupARN>arn:aws:autoscaling:us-east-1:123456789012:autoScalingGroup:66
          be2dec-ee0f-4178-8a3a-e13d91c4eba9:autoScalingGroupName/my-test-asg<
      </AutoScalingGroupARN>
        <TerminationPolicies>
          <member>Default</member>
        </TerminationPolicies>
        <MaxSize>5</MaxSize>
      </member>
    </AutoScalingGroups>
  </DescribeAutoScalingGroupsResult>
  <ResponseMetadata>
    <RequestId>cb35382a-64ef-11e2-a7f1-9f203EXAMPLE</RequestId>
  </ResponseMetadata>
</DescribeAutoScalingGroupsResponse>

```

List of Tables

2.1	Comparison of load balancers	14
2.2	Comparison of hypervisors resizing capabilities	15
2.3	Comparison of cloud providers scaling capabilities	16
3.1	Comparison of cloud providers approach to adaptivity	20
5.1	Monitoring - summary of issues	29
7.1	Configuration of hardware/virtual machines used during tests	33
7.2	Price for a stack in the given cloud provider	34
7.3	Chosen cloud providers for the given stack	35
7.4	Configuration of hardware/virtual machines used during tests	40
7.5	Average deployment time for the service with the various number of VMs used for the whole environment	51

List of Figures

1.1	Public Cloud Services Market Size, 2010-2016 (forecast). Source: <i>Gartner</i> , 08/2012	8
2.1	Scalability layers	11
2.2	Amdahl's law	12
3.1	Elasticity controller	17
3.2	Threshold model	18
4.1	Major obstacles to cloud adoption in 2013 according to [16]	22
4.2	Top benefits of using the hybrid model according to [17]	23
5.1	Cloud-SAP high level overview	26
5.2	Layered structure of Cloud-SAP	26
5.3	autonomic component	27
5.4	Cloud-SAP as an hierarchical autonomic system	28
7.1	Deployment cost: environment configuration	34
7.2	Comparison of the deployment cost when the service is deployed only on a selected cloud provider or a combination of cloud providers selected by Cloud-SAP	35
7.3	Auto-scaling - single-provider based: CPU usage function	37
7.4	Auto-scaling - single-provider based: deployment diagram of Carina	39
7.5	Auto-scaling - single-provider based: deployment diagram of Cloud-SAP	39
7.6	Auto-scaling - single-provider based: environment configuration	40
7.7	Auto-scaling - single-provider based: comparison of cost/CPU usage	42
7.8	Auto-scaling - multiple-provider based: CPU usage function	44
7.9	Auto-scaling - multiple-provider based: deployment diagram of Carina	45
7.10	Auto-scaling - multiple-provider based: deployment diagram of Carina	46
7.11	Auto-scaling - multiple-provider based: environment configuration	47
7.12	Auto-scaling - multiple-provider based: comparison of processed transaction per second	48
7.13	Deployment time - solution comparison: physical environment setup	50
7.14	Deployment diagram of <i>Carina</i>	50
7.15	Deployment diagram of <i>Cloud-SAP</i>	51
7.16	Average deployment time for two competing products when the variable is the number of instances of VMs	52

Bibliography

- [1] Apache httpd. <http://httpd.apache.org/>. Accessed: 2013-08-30.
- [2] Apache mod_jk. <http://tomcat.apache.org/download-connectors.cgi/>. Accessed: 2013-08-30.
- [3] Apache mod_proxy_balancer. http://httpd.apache.org/docs/2.2/mod/mod_proxy_balancer.html. Accessed: 2013-08-30.
- [4] Auto scaling - amazon web services. <http://aws.amazon.com/autoscaling/>. Accessed: 2013-08-30.
- [5] Cloudstack. <http://cloudstack.apache.org/>. Accessed: 2013-08-30.
- [6] Cloudwatch. <http://aws.amazon.com/es/cloudwatch/>. Accessed: 2013-08-30.
- [7] Eucalyptus. <http://www.eucalyptus.com/>. Accessed: 2013-08-30.
- [8] F5. <http://www.f5.com/>. Accessed: 2013-08-30.
- [9] Haproxy. <http://haproxy.1wt.eu/>. Accessed: 2013-08-30.
- [10] Haproxydoc. <http://haproxy.1wt.eu/download/1.3/doc/configuration.txt>. Accessed: 2013-08-30.
- [11] Opennebula carina. <https://github.com/blackberry/OpenNebula-Carina>. Accessed: 2013-12-16.
- [12] Opennebula project. <http://opennebula.org/>. Accessed: 2013-08-30.
- [13] Openshift. <https://www.openshift.com/>. Accessed: 2013-08-30.
- [14] Openshiftscaling. <https://www.openshift.com/developers/scaling>. Accessed: 2013-08-30.
- [15] Openstack. <http://www.openstack.org/>. Accessed: 2013-08-30.
- [16] The future of cloud computing 3rd annual survey. Tech. rep., 2013. Report from the survey available at <http://mjskok.com/resource/2013-future-cloud-computing-3rd-annual-survey-results>. Accessed: 2013-09-03.
- [17] Rackspace 2013 hybrid cloud survey results. Tech. rep., 2013. Report from the survey available at http://www.rackspace.com/knowledge_center/article/rackspace-2013-hybrid-cloud-survey-results. Accessed: 2013-09-03.
- [18] A., B. Characteristics of scalability and their impact on performance.

- [19] ABDEEN, M., AND WOODSIDE, C. M. Seeking optimal policies for adaptive distributed computer systems with multiple controls. Third International Conference on Parallel and Distributed Computing, Applications and Technologies.
- [20] ABDELZAHER, T., SHIN, K., AND BHATTI, N. Performance guarantees for web server end-systems: A control theoretical approach, 2002.
- [21] BUYYA, R., RANJAN, R., AND CALHEIROS, R. Intercloud: Scaling of applications across multiple cloud computing environments.
- [22] BUYYA, R., RANJAN, R., AND RN., C. Intercloud: Utility-oriented federation of cloud computing environments for scaling of application services.
- [23] BUYYA, R., YEO, C., AND VENUGOPAL, S. Market-oriented cloud computing: Vision, hype, and reality for delivering it service s as computing utilities. *High Performance Computing* (2008).
- [24] BUYYAA, R., YEOA, C., VENUGOPALA, S., BROBERGA, J., AND BRANDICC, I. Cloud computing and emerging it platforms: Vision, hype, and reality for delivering computing as the 5th utility. *Future Generation Computer Systems* 25 (2009).
- [25] DIAO, Y., LUI, X., FROEHLICH, S., HELLERSTEIN, J., AND PAREKH, S. On-line response time optimization of an apache web server. International Workshop on Quality of Service.
- [26] DUBOC, L., ROSENBLUM, D., AND WICKS, T. A framework for modelling and analysis of software systems scalability. Proceeding of the 28th international conference on Software engineering.
- [27] G., A. Validity of the single processor approach to achieving large scale computing capabilities. vol. 983 of *AFIPS Conference Proceedings*.
- [28] HILL, M. What is scalability? ACM SIGARCH Computer Architecture News.
- [29] IBM. An architectural blueprint for autonomic computing, 4th edition.
- [30] JIANG, Y., PERNG†, C., LI, T., AND CHANG, R. Asap: A self-adaptive prediction system for instant cloud resource demand provisioning. IEEE International Conference on Data Mining.
- [31] LEAVITT, N. Is cloud computing really ready for prime time? Available at <http://www.leavcom.com/pdf/Cloudcomputing.pdf>. Last accessed: 2013-09-10.
- [32] LEONG, L., TOOMBS, D., GILL, B., PETRI, G., AND HAYNES, T. Magic quadrant for cloud infrastructure as a service. Tech. rep., Gartner, Gaithersburg, Aug. 2013.
- [33] LITOIU, M., WOODSIDE, M., AND ZHENG, T. Hierarchical model-based autonomic control of software systems. Design and evolution of autonomic application software.
- [34] MAO, M., LI, J., AND HUMPHREY, M. Cloud auto-scaling with deadline and budget constraints. IEEE/ACM International Conference on Grid Computing.
- [35] MASSIE, M., CHUN, B., AND CULLER, D. The ganglia distributed monitoring system: design, implementation, and experience.
- [36] MELL, P., AND GRANCE, T. The nist definition of cloud computing.
- [37] R., M. Autonomic computing.

-
- [38] SHEN, Z., SUBBIAH, S., GU, X., AND WILKES, J. Cloudscale: Elastic resource scaling for multi-tenant cloud systems. ACM Symposium on Cloud Computing – SOCC2011.
 - [39] STECCA, M., BAZZUCCO, L., AND MARESCA, M. Sticky session support in auto scaling iaas systems. IEEE World Congress on Services.
 - [40] VAQUERO, L., RODERO-MERINO, L., AND BUYYA, R. Dynamically scaling applications in the cloud, 2011.
 - [41] ZHENG, T., YANG, J., AND WOODSIDE, M. Tracking time-varying parameters in software systems with extended kalman filters.