

A Generalizable Framework for Automated Cloud Configuration Selection

Supervisors: Adam Barker & Yuhui Lin

Jack Briggs - 140011358

MSc Data-Intensive Analysis

2019-06-06

Abstract

Outline of the project using at most 250 words

Declaration

I declare that the material submitted for assessment is my own work except where credit is explicitly given to others by citation or acknowledgement. This work was performed during the current academic year except where otherwise stated. The main text of this project report is NN,NNN* words long, including project specification and plan. In submitting this project report to the University of St Andrews, I give permission for it to be made available for use in accordance with the regulations of the University Library. I also give permission for the title and abstract to be published and for copies of the report to be made and supplied at cost to any bona fide library or research worker, and to be made available on the World Wide Web. I retain the copyright in this work.

Contents

1	Introduction	1
1.1	Cloud Computing	1
1.1.1	Cloud optimization	2
1.1.2	Benefits	2
1.1.3	Challenges	2
1.2	Aims and Objectives	5
1.3	Contributions	6
1.4	Dissertation Outline	6
2	Literature Survey	7
2.1	Cloud services	7
2.2	Cloud variability	7
2.3	Optimization methods	7
2.3.1	CherryPick	7
2.3.2	PARIS	7
2.3.3	Ernest	7
2.3.4	Daleel	7
2.3.5	Okta	7
2.3.6	OTHERS, LOOK UP	7
2.3.7	Exhaustive search	8
2.4	Benchmarks	8
2.4.1	Cloudsuite	8
2.4.2	vBench	8
2.4.3	OTHERS, LOOK UP	8

2.5	Infrastructure-as-code	8
2.5.1	Terraform	8
2.5.2	Apache Libcloud	8
2.5.3	Chef	8
2.5.4	Puppet	8
3	Requirements specification	9
3.1	Use-case	9
3.1.1	Bayesian Optimization	9
3.1.2	Video transcoding	9
3.2	Requirements	9
3.3	Optional Requirements	9
4	Generic Automated Cloud Configuration Optimization Framework	10
4.1	Motivations	10
4.1.1	Numerical Optimization	10
4.1.2	Modularity	11
4.2	System Architecture	11
4.3	Inputs and outputs	13
4.4	Searcher	14
4.5	Selector	16
4.5.1	Exact vs. Closest Match	16
4.5.2	Encoding the Search space	17
4.6	Deployer	18
4.6.1	VM Provision	18
4.6.2	Docker Deployment	18
4.6.3	Ping server	19
4.6.4	Simulated Deployment	19
4.6.5	Interpreter	19
4.7	Interpreter	19
5	Bayesian VM Optimization System	20
5.1	Bayesian Optimization	20
5.2	Benchmark Deployment	20

5.3	Search space encoding	20
5.4	Objective Measure	20
6	Implementation	21
6.1	General usage	21
6.2	Driver	21
6.3	Searcher	21
6.3.1	Spearmint	21
6.4	Selector	21
6.4.1	Exact Match	21
6.5	Deployer	21
6.5.1	VM Provisioner	22
6.5.2	Terraform	22
6.5.3	Docker deployer	22
6.5.4	Cloudsuite	22
6.5.5	Ping servers	22
6.6	Interpreter	22
6.6.1	Sysbench	22
6.6.2	Cloudsuite	22
6.6.3	vBench	22
6.6.4	Fake Deploy	22
7	Evaluation	23
7.1	Evaluation Approach	23
7.1.1	Framework evaluation	23
7.1.2	Bayesian Optimization	23
7.1.3	Evalutaion objectives	23
7.2	Results Analysis	23
7.2.1	Exhaustive search	23
7.2.2	Bayesian Optimization	23
8	Related and Future work	24
9	Conclusion	25

List of Figures

4.1	A flowchart showing the processes involved in and information flow through the designed system.	12
4.2	An example of BO's working process. Taken from Figure 5 in [1]	16

Chapter 1

Introduction

1.1 Cloud Computing

Cloud computing is an ever-growing field that now ranges from Infrastructure-as-a-service (IaaS) to Software-as-a-service(SaaS). Services under cloud computing are characterised by their ability to offer access to a shared pool of highly elastic on-demand computing resources that offer broad network access [2,3]. Cloud services as an industry has had an explosive growth, and it has been predicted that 83% of enterprise (Companies with 1000+ employees) workloads will be in the cloud by 2020 [4], with 41% run on public cloud platforms such as Amazon AWS and Microsoft Azure. Services offered range from various levels and forms of abstractions, from directly provisioning Virtual Machines (VMs) or storage services, allowing users full control over their cloud infrastructure, to deploying 'serverless' containers, where the actual managing of the hardware is instead handled by the cloud provider.

The appeal is obvious, with cloud services allowing organizations and developers to utilize a diverse range of computational resources on demand without any up-front commitment or cost [5]. This can lead to both significant cost-savings as well as improved revenue through better customer experiences and enabling risk-free experimentation [6]. Academics, too, are utilizing the available services as volumes of data grow impractical to store and analyse on local machines [7,8]. This includes large-scale collaborative projects involving huge data sets hosted on the cloud such as the 1000 Genomes Project¹ or Common Crawl².

¹<https://aws.amazon.com/1000genomes>

²commoncrawl.org

1.1.1 Cloud optimization

A wide range of applications are now deployed on cloud machines or make use of objects stored on them, from large-scale data analytics jobs mentioned to media-streaming servers such as Netflix or Twitch [9]. The resource dependencies of these applications similarly vary widely, from the CPU dependent data analysis tasks to network-heavy streaming services. Virtual machines offered by different cloud providers vary in terms of memory amounts of number and speed of virtual CPUs (vCPUs), and each application's performance will have different relationships with these options. While medium-length video transcoding operations will benefit primarily from faster processing speeds, data analysis tasks involving large datasets may find a more cost-effective option in prioritising VMs with a local solid-state drive (SSD) offering high I/O performance. Non-critical batch workloads can often benefit from using 'pre-emptible' or 'spot' instances, which offer large discounts on the condition that your machine can be terminated with little notice to free up resources.

1.1.2 Benefits

It is desirable for both users and providers to maximise the optimize purchased cloud configurations to best serve the needs of their applications. Users or developers who fail to do this risk paying far more than they need to for the same performance. A given data analysis task can cost around 3.4 times as much on an average configuration compared to the optimal available option [1]. Even serverless frameworks simply shift the burden of optimization from the users to the cloud providers. For cloud providers too, efficient deployment across available Virtual Machines frees up extra resources available for other purposes or other customers. Alternatively, identifying and avoiding the co-location of CPU intensive workloads can reduce resource contention amongst users, leading to improved performance [10]. In addition, energy-related costs make up to 42% of managing a data-centre, and the ability to idle inactive resources would lead to a significant reduction both in energy cost and environmental impacts. [11,12].

1.1.3 Challenges

Picking an optimal cloud configuration for any given application is far from trivial. 'Larger' VM types may not provide improved performance [13], or may do so at a far less cost-efficient rate. Even if one knows their exact budget along with the constraints placed on their application's performance, they must still find a way to search or model across the entire range of cloud providers to ensure they have found the optimal cloud configuration for their application. This is further complicated by the variety of instance types available across multiple providers, and the variation in performance between virtual machines of

the same specification due to differences in underlying hardware, network traffic, and co-located tenants located on the same physical machine.

Search space

The diversity in services provided by different cloud providers creates a large search space. At the time of writing, Amazon EC2³ alone offers over 200 predefined instance models⁴, while Google Compute Engine⁵ allows users to define their own machine types, ranging from small VMs with 1 vCPU and 1 GB of RAM to 96 vCPUs with 624 GB of RAM⁶, and includes options for specifying the CPU platform or adding GPUs. VMs can also differ in terms of network performance, local storage, and CPU speed, underlying hardware, and can start able to run a variety of operating systems and associated software.

The search space is not just challenger in terms of range, but is also hard to define in terms of constraints. Many hardware options are precluded by or only permitted alongside other options. To simplify selection, leading providers generally categorize instance types are generally split into categories such as 'compute-optimized' and 'memory-optimized' machines, each of which is made up of machines of various 'sizes,' such as 'tiny,' 'small,' 'large,' referring to the computational resources available to that virtual machine. Even with this categorization, it can be hard to formalize or encode the search space in such a way to best serve an automated optimization software.

Hardware and Software Heterogeneity

Despite two Virtual Machines sharing the exact same configuration, if their underlying machines possess different physical hardware then their performance can differ significantly [14]. Amazon's M4 instance types, for example, can come with either 2.3 GHz Intel Xeon® E5-2686 v4 (Broadwell) or a 2.4 GHz Intel Xeon® E5-2676 v3 (Haswell), based on the machine they happen to be hosted by⁷. However, more recent studies involving newer instance types have shown a dramatic reduction in this CPU-based variability for certain instance types [15, 16].

Problems regarding hardware heterogeneity expand further when searching across multiple providers. Different providers may use vary different hardware for machines reporting the same number of vCPUs or disk speeds. Software too, can differ. Amazon instances come configured according to a template referred

³<https://aws.amazon.com/ec2/>

⁴<https://www.ec2instances.info>

⁵<https://cloud.google.com/compute/>

⁶<https://cloud.google.com/custom-machine-types/>

⁷<https://aws.amazon.com/ec2/instance-types/>

to as an Amazon Machine Image (AMI)⁸, while Google Compute Engine instances use a predefined boot disk. These images include an operating system and associated software, storage snapshots, launch permissions, and device mappings. The APIs used to communicate with cloud providers in order to set up cloud services in the first place similarly differ between providers. Even when using Infrastructure as code (IaC) tools, the nomenclature and configuration file structure can differ widely.

When searching across multiple instance types, the hardware and software heterogeneity must be taken into account, especially when the search spans separate cloud providers.

Variation

Even when run on identical cloud configurations from the same provider, the same application or benchmark can have significantly different performance, more so than is observed when run in on-premise environments [14]. This is in part due to the hardware heterogeneity in underlying machines described above, but can also occur due to multitenancy. Multiple customers may share computational resources by running VMs on the same physical machine. This can lead to a slowdown in one instance due to the behaviour of another colocated tenant. This problem is referred to as the 'noisy neighbour' problem [12]. Disk I/O operations or network traffic between instances and storage disks could similarly be affected by the machine's relative locations within or between datacentres.

All the above problems causing variation within instance types are generally out of the control of the customer. Cloud providers may develop methods of minimising the problems, shown by the reduction in the variation of intra-cloud network performance [17] and CPU performance between identical instance types [15, 16]. Nonetheless, in cloud configurations where this variation is still significant, it adds a significant element of randomness to the performance of various applications on cloud services. This randomness must be accounted for by any method used to optimize cloud configurations, and rules out assumptions of deterministic application performance.

Application range

As mentioned earlier, numerous types of applications are deployed on the cloud, from long-running web-servers and media-streaming services to one-time data analytics jobs. The performance and cost-effectiveness of these applications can have very different non-linear relationships with cloud configurations, making them difficult to model analytically, even if the contents of the application are known

⁸<https://docs.aws.amazon.com/AWSEC2/latest/UserGuide/AMIs.html>

exactly. In some cases, such as in serverless container services such as Google Cloud Run⁹ the application is not known at all, and must be treated as a black box.

To ensure an optimization method is as generalizable as possible, it must be able to utilize user's specifications and the outputs of an application, but must treat the application itself as a black box, and not make hard assumptions regarding any of its characteristics or relationships with configuration options.

Objective measure

Finally, the very meaning of the word optimal can depend on the application and the user. The ultimate goal is far from fixed, as optimizing for a good mean performance for a server may fail to optimize for high stress situations, leading to lower-than expected performance under these conditions. Even in single batch jobs, different users may tolerate different reductions in performance for different cost-savings. While an effective option is to simply present to the user an informative representation of how cost and performance vary with cloud configuration options, this is not possible in many black-box options, and limits the degree of automation in the optimization process.

To truly automate an optimization process for any provided application, an 'objective measure' must be provided, dependent on cost and performance, which reflects the ultimate goal of the optimization process. An optimization method must be able to accommodate different user requirements through flexibility regarding the objective measure used.

1.2 Aims and Objectives

The aim of this project is both to present a framework for completely automated selection of an optimal cloud configuration for any given application, as well as to deliver a fully functional implementation of this framework for an indicative use-case. To achieve this aim, the framework must achieve the following objectives:

- Present a fully automated process for optimization cloud configuration.
- Be applicable for any form of cloud application.
- Not be limited to any one cloud provider or service.
- Take into account variation within and between instance types.

⁹<https://cloud.google.com/run/>

The framework should require only the following inputs, from which it should output a single estimate for the optimal cloud configuration:

- The application itself, to be set up on a given cloud service.
- Details of the search space and its constraints.
- An objective measure with which to evaluate cloud configurations.

In addition, we wish to deliver a fully functional implementation for an indicative use case, which will achieve the following objectives:

- Automate selection of cloud configuration samples based on a search space of cloud configurations spanning multiple providers.
- Provision specified instance types from multiple providers based on a given cloud configuration.
- Deploy and return logs for a given application on a given cloud service based on a given cloud configuration.
- Interpret logs based on the performance and cost of a given application on a given cloud configuration.

1.3 Contributions

1.4 Dissertation Outline

Chapter 2

Literature Survey

2.1 Cloud services

2.2 Cloud variability

2.3 Optimization methods

2.3.1 CherryPick

2.3.2 PARIS

2.3.3 Ernest

2.3.4 Daleel

2.3.5 Okta

2.3.6 OTHERS, LOOK UP

asdf

2.3.7 Exhaustive search

2.4 Benchmarks

2.4.1 Cloudsuite

2.4.2 vBench

2.4.3 OTHERS, LOOK UP

2.5 Infrastructure-as-code

2.5.1 Terraform

2.5.2 Apache Libcloud

2.5.3 Chef

2.5.4 Puppet

Chapter 3

Requirements specification

3.1 Use-case

3.1.1 Bayesian Optimization

3.1.2 Video transcoding

3.2 Requirements

3.3 Optional Requirements

Chapter 4

Generic Automated Cloud Configuration Optimization Framework

4.1 Motivations

4.1.1 Numerical Optimization

In any optimisation process, the value of some objective function is minimised or maximised by adjusting its input variables or parameters. An optimisation algorithm begin with an initial guess of these variables and iterates through improved estimates until they terminate, hopefully providing an estimated solution [18]. An optimization function effectively takes previous an objective function, along with its previous inputs and outputs as arguments, and outputs either the next sample to take, or the predicted best input if the stopping condition is met.

In our generalized case, the optimisation method and objective function are both unknown, but we know that our objective function will always involve selecting some cloud configuration based on the inputs that describe that configuration, deploying some application onto this configuration, and interpreting its performance to give some objective measure.

4.1.2 Modularity

In our generalized case, we do not yet know the optimization method or objective function. As shown in our Literature review, a number of optimization algorithms have been attempted and are available, such as exhaustive search, Bayesian optimization [1], analytical methods that directly model the application’s relationships with configuration variables [19], and a data-driven approach predicting performance by its relationship with pre-run benchmarks [13].

Much like how our objective measure will differ from application to application, the best optimization method may also differ depending on the type of application tested. Applications whose performance is accurately reflected by known benchmarks may benefit from PARIS’s data-driven approach, while situations with unclear non-linear relationships to configuration variables will be well suited to Bayesian optimization. Because of this, while we will focus on developing an implementation of Bayesian optimization, as per our use case, it is desirable to ensure that our example optimization method can be replaced by another if the user deems it more appropriate to their use-case.

Because of this, we are left with multiple parts of our system which are highly dependent on a user’s individual use-case, namely the objective function, the objective measure to be optimized, the search space constraints, and the application itself. Because of this, we have taken a highly modular approach to our design. While our implementation should be self-sufficient to run its example case, it should also be made up of several replaceable components for future extensions.

4.2 System Architecture

Our optimization process can be broken down into four replaceable components: A Searcher, which runs the optimisation algorithm, testing out various inputs in an attempt to maximise or minimise the objective function; a Selector, which interprets the inputs to determine what cloud configuration is being tested; a Deployer, which provisions the machines needed for that cloud configuration, deploys the application, and once it has terminated returns any required logs from it; and an Interpreter, which takes these logs to calculate the objective measure which is returned to the Searcher as the returned value for the objective function. A diagram of this breakdown is shown in figure 4.1. The components themselves are examined in further detail in the following sections.

Overall, the system performs an optimization algorithm shown in Procedure 1, such as Bayesian optimization, coordinate descent, random search, or exhaustive search, and drives the optimization process by iterating through potential input variables. For each set of these inputs, it take a sample from a single 'job,' where it runs through a single loop of the Selector, Deployer, and Interpreter components. The

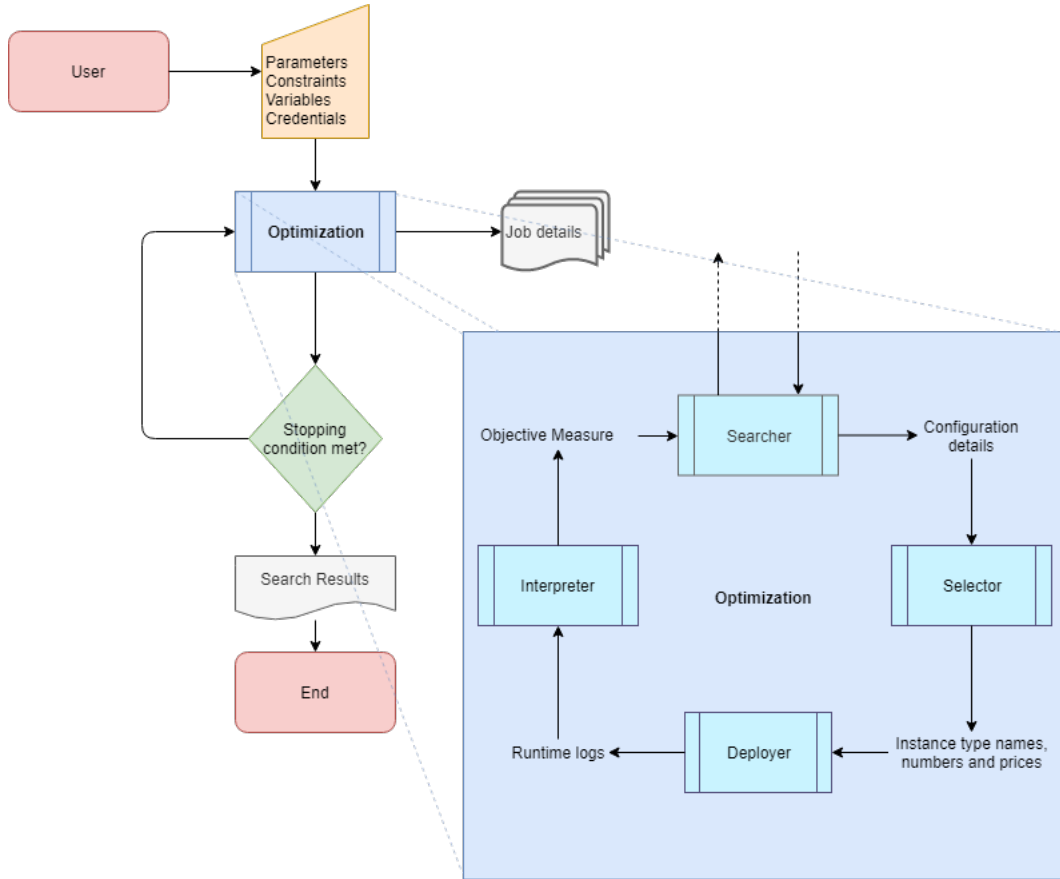


Figure 4.1: A flowchart showing the processes involved in and information flow through the designed system.

Searcher component then uses the results from inputs tried so far to choose the next set of input variables to sample, as well as provide the current estimate of the best possible set of input variables to minimize or maximize that result. Once the stopping conditions are met, the latest estimate for the best input variables is output by the system.

Algorithm 1 Optimization Procedure

```

procedure OPTIMIZATION(Starting variables, Stopping conditions, Searcher, Selector,
                        Deployer, Interpreter)
   $\vec{x}_0 \leftarrow \text{Starting variables}$ 
  for  $i = 0, 1, 2, 3, \dots$  do
     $Config_i, Price_i \leftarrow Selector(x_i)$ 
     $logs_i \leftarrow Deployer(x_i, Config_i)$ 
     $results_i \leftarrow Interpreter(x_i, logs_i, Price_i)$ 
     $\vec{x}_{i+1}, stop, best\vec{x} \leftarrow Searcher((x_-, x_1, \dots, x_i), (results_0, results_1, \dots, results_i),$ 
                                          Stopping conditions)

    if stop then
      return  $best\vec{x}$ 
    end if
  end for
end procedure

```

It is assumed that in the vast majority of cases, the user would provide their own Interpreter and Selector. Interpreters and Selectors are very dependent on the form these logs will take, and the form the search space will take, and are extremely hard to generalise. For this reason, the modular design of our solution should make it simple for any component to be supplied or replaced by the user. In only some cases should the user need to provide their own Deployer. Applications can often be contained within Docker containers, and aside from occasional setup, for example in multi-node clusters or multi-container applications, a Deployer which provisions a given configuration from a given provider, and then deploys and attaches to a user-provided docker image to collect its logs will be sufficient. Only in rare cases would the user be required to provide their own Searcher. This is because optimization algorithms can be applied to any deployment, with only small modifications necessary in rare cases for specific cases, and an implementation of Bayesian Optimization will be provided.

4.3 Inputs and outputs

The user inputs are shown in figure 4.1 as being made up of 4 parts: Parameters, constraints, variables, and credentials. The parameters would be options for the searching method itself, such as the stopping conditions or maximum number of concurrent jobs. Variables and constraints together define the possible search space for the input variables, as well as their starting values. Each variable must be given a type,

such as string or integer, and a set of constraints such as maximum or minimum for numeric variables or the available options for categorical variables such as strings. Finally, the user would have to submit paths to the files holding the credentials to allow the program to provision and access virtual machines on the tested cloud providers. It is of course essential that credentials are never stored in any form by the program itself, and only paths to the relevant files are passed to the relevant tools and interfaces.

The figure also shows the relevant outputs created by the system, in the form of 'job details' and 'search results' output as files containing details describing the jobs performed and their individual results, as well as the final prediction for the optimal configuration produced by the system. Keeping outputs from individual jobs is important, both for debugging the system during development or component switch-over, as well as to allow users to avoid repeating unnecessary repeats of already tested configurations when performing multiple search experiments.

Algorithm 1 shows the general form the optimization process takes with only one concurrent job allowed at any time, showing the expected input and output of each individual component. These are explained in more detail within each component's individual section below.

4.4 Searcher

The Searcher component is actually responsible for handling the optimization algorithm itself. The searcher takes the results from previous jobs and determines either the next best sample to take within the available search space, or to terminate the process and return the current best prediction for the optimal configuration. A perfect Searcher is one that returns a perfect prediction of the optimal configuration in as few samples as possible. Fewer samples means a reduction in the time and cost taken to perform the search. Often, there is a trade-off here, where extending a search to include more and more samples, such as by setting less restrictive stopping conditions, increases its predictive accuracy. Some search methods, such as PARIS [13], can leverage past search experiments to reduce the number of samples taken. Generally however, the stopping conditions and search method will control this trade-off, and can be set according to the user's search budget or time constraints.

In order to accurately predict the best next sample to take, the Searcher must know the available search space. A description of how to encode cloud configurations into a set of variables and constraints is done in the Selector section. Ideally, we want our searcher to be capable of performing multiple jobs concurrently.

Here we describe and compare the design considerations of two important examples of Searchers which will be used for the evaluation of our implementation.

Exhaustive search

In an exhaustive search, every possible combination of the inputs is sampled, giving a complete analysis of the entire search space. This obviously takes many samples, $n * \prod_{i=1}^J x_i$ where x_i is the number of options for the i th of J variables, and n is the number of samples taken from each configuration. This results in a large or even infinite search cost and time, but is almost certain to return the optimal result, depending on the amount of randomness involved in sampling.

Bayesian Optimization

Bayesian Optimization is an optimization method specifically designed for situations where the objective functions itself is unknown beforehand and expensive to perform, but can have its outputs directly observed through experiments [20]. It models the objective function as a stochastic process, a 'prior function,' and computes its confidence intervals based on samples. Using these confidence intervals and a pre-defined acquisition function it decides future samples to take based on where there it estimates there to be the most potential gain. A diagram showing this process is shown in figure 4.2. By this process Bayesian Optimization can find optimal or near-optimal solutions for any non-parametric problem in a relatively small number of samples compared to other optimization methods. In addition, whereas other methods, such as exhaustive search, may handle uncertainty through sampling results from the same inputs multiple times, Bayesian Optimization can incorporate this uncertainty into its estimates, further reducing the number of samples needed.

There are a number of possible prior functions, acquisition functions, and stopping conditions that can be used with BO, and the Cherrypick paper goes into detail on the reasoning behind which options are best for cloud configuration selection specifically [1]. Some notable differences in our case, however, is that CherryPick was specifically focused on batch jobs, where what is measured is simply a function of an instance type's cost and its time taken to perform a given job. Its acquisition function is specified to this purpose, minimizing costs but biased towards configurations which satisfy a soft performance constraint by successfully performing the batch job within a given time. In our case, our acquisition function must be more general, and we will therefore be relying on the user to ensure that whatever objective measure it returned by their Interpreter has already taken into account soft or hard performance constraints such as this.

4.5 Selector

The Selector interprets the variables provided by the Searcher component into the form of an available cloud configuration. Cloud configurations have a number of variables that can describe them, such as vCPU number, memory amount, disk speed, number of instances, instance category, machine type, and cloud provider. The selector must use whatever combination of these is provided and find either the exact or most similar cloud configuration available, passing this information on to the Deployer. In the tools used for our implementation, no cross-provider API was available to directly translate machine specifications into the virtual machine types available from different providers.

4.5.1 Exact vs. Closest Match

For finding an exact match, the Selector can simply lookup the appropriate instance type from a dataset according to the input variables stored as each instance type’s attributes. Looking for a closest match rather than an exact one gives more flexibility in how the input variables can be encoded, but means more complicated decisions such as attribute priority must be made, and extra assurances made to not repeat unnecessary samples when multiple sets of inputs describe the same closest input type. For our use-case, we will specifically use the simpler Exact Match method, but thanks to the modular design it will be easy

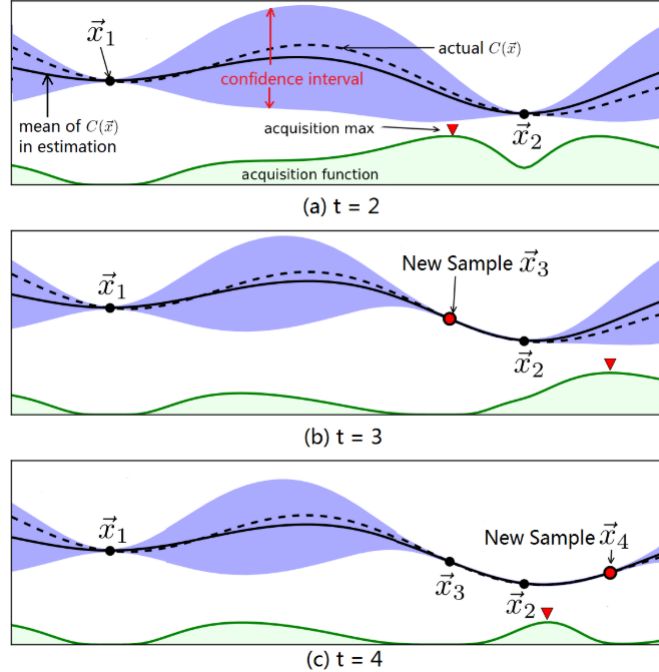


Figure 4.2: An example of BO’s working process. Taken from Figure 5 in [1]

Instance Type	Provider	Category	vCPU Number	Memory
n1-standard-2	GCE	General	2	7.50
n1-standard-4	GCE	General	4	15.00
n1-standard-8	GCE	General	8	30.00
c5.large	EC2	CPU	2	4.00
c5.xlarge	EC2	CPU	4	8.00
c5.2xlarge	EC2	CPU	8	16.00

Table 4.1: A possible way of separating instance types into 4 descriptive variables. Providers were either the Google Compute Engine (GCE) or the Amazon Elastic Compute Cloud (EC2)

to extend to include a Closest Match Selector option if time allows.

4.5.2 Encoding the Search space

Whether using exact or closest match, it must be decided how to encode cloud configurations into a set of input variables. This problem is further complicated by the fact that constraints for certain inputs may depend on the values of others, as is often the case in cloud computing. Large memory amounts are often only available on machines with more vCPUs, and some providers may offer available configurations others do not. Google Cloud Platform allows users to specify custom machine types, but even these do not allow any possible combination (for example, Memory constraints are tied to vCPU number, and vCPU number must be divisible by 2).

Despite these problems, there are at least clear patterns prevailing throughout leading cloud providers, such as separating machine types into categories equivalent to 'Compute-optimised', 'Memory-optimised', and 'Storage-optimized,' each with a set of machines with between 2 and 96 CPUs. In lieu of a cross-provider service to match a given specification to a specific cloud instance type, something which would be outside the scope of this project, these industry-standard categorisations can be used to encode the search space in a reasonable manner. Table 4.1 shows an example of how we have used these patterns to encode 6 instance types into a set of 3 easily interpreted variables; Provider, vCPU number, memory, and machine category. A searcher tool could easily filter a dataset of this form to find the instance-type for a set of input variables.

In the end, however, the important features and constraints for the search space will differ for each user, and it may be beneficial to run multiple experiments in different search spaces before settling on a final decision for an instance type. While we provide in our associated implementation an example of how to encode and select available cloud configurations for our Bayesian Optimization tool, we think it best to ultimately leave it to the user to design and implement a Selector system that works well for their specific use-case.

4.6 Deployer

The Deployer deploys the user-provided application, batch job, or benchmark onto the selected cloud configuration, and collect any necessary analysis from it. Typically this will involve provisioning the necessary machines from the given provider, followed by deploying the given application onto these machines, and either collecting logs from them or from a networked instance or cluster.

4.6.1 VM Provision

Aside from in serverless computing, the first step of deploying any cloud-based application is likely to be provisioning the virtual machines themselves from a cloud provider. The Deployer should be capable of requesting any virtual machine chosen by the Selector, regardless of provider. This can be simplified using Infrastructure as code (IaC) tools such as Terraform, which offers the ability to codify the APIs from many different providers into declarative configuration files. As long as configuration files and credentials for use by an IaC tool are supplied for each possible provider, then a Deployer can call the corresponding IaC tool to provision the machines. There are several requirements for the IaC tool that is used. As instance type and number of machines will be supplied by the Selector, and therefore cannot be known until the machines are provisioned, either they must be declarable as variables when the tool is run, or the configuration files must be suitable for automated editing just beforehand. Along with this, the tool must support multiple concurrent tasks with their own specifications and outputs, if we hope to allow the Searcher to take multiple samples at once.

4.6.2 Docker Deployment

Once the virtual machines themselves have been provisioned, the application must be deployed onto them. In the most basic case, we assume that the application being tested is available in the form of a Docker image, and that a single instance is provisioned on which to run that application. and that the machines provisioned operate as either a Kubernetes based cluster or a single instance. While other forms of application or cluster architecture may be used, ready-made Kubernetes clusters are available on several major cloud providers. The modular design allows users to implement their own Deployers to deal with alternative situations.

Both Kubernetes and Docker offer remote APIs. For Kubernetes clusters available on cloud providers, no set up is required, while for single machines the VM provisioner must install Docker and direct it to a public-facing port. The remote APIs can be utilized as long as the VM Provisioner is capable of

returning the public facing IP of the provisioned machines, and that security credentials are provided. The Deployer must attach to any deployment, or wait for its completion, and obtain and return any and all logs it produces for use by the Interpreter.

4.6.3 Ping server

4.6.4 Simulated Deployment

It may be advantageous, especially during debugging or if using a 'closest match' approach to configuration selection, to instead either simulate the response from provisioning and deployment of an application, or to return a previously sampled response. To this end, it may be useful to ensure full logs of any deployment are stored so that, if desired, future calls to 'Deploy' an already sampled configuration for a given application can simply be responded with a randomly distributed value based on previous results from the same configuration.

4.6.5 Interpreter

The Interpreter must interpret whatever logs and other information is returned by the Deployer, along with the cost of the cloud configuration provided by the Selector, in order to return an objective measure for the sampled cloud configuration. It is this returned value that the Searcher will be attempting to minimise or maximize. We leave it to the user to develop an Interpreter for their use-case. This will likely involve extracting relevant information from the deployed application's logs, and applying constraints based on the time taken or machine pricing.

4.7 Interpreter

Chapter 5

Bayesian VM Optimization System

5.1 Bayesian Optimization

5.2 Benchmark Deployment

5.3 Search space encoding

5.4 Objective Measure

Chapter 6

Implementation

6.1 General usage

6.2 Driver

6.3 Searcher

6.3.1 Spearmint

6.4 Selector

6.4.1 Exact Match

6.5 Deployer

asdf

6.5.1 VM Provisioner

6.5.2 Terraform

6.5.3 Docker deployer

6.5.4 Cloudsuite

6.5.5 Ping servers

6.6 Interpreter

6.6.1 Sysbench

6.6.2 Cloudsuite

6.6.3 vBench

6.6.4 Fake Deploy

Chapter 7

Evaluation

7.1 Evaluation Approach

7.1.1 Framework evaluation

7.1.2 Bayesian Optimization

7.1.3 Evalutaion objectives

7.2 Results Analysis

7.2.1 Exhaustive search

7.2.2 Bayesian Optimization

Cross-provider

Concurrent Jobs

Chapter 8

Related and Future work

Chapter 9

Conclusion

[21]

Bibliography

- [1] O. Alipourfard, H. H. Liu, J. Chen, S. Venkataraman, M. Yu, M. Zhang, Y. University, H. Harry Liu, J. Chen, S. Venkataraman, M. Yu, and M. Zhang, “CherryPick: Adaptively Unearthing the Best Cloud Configurations for Big Data Analytics,” in *Proceedings of the 14th USENIX Conference on Networked Systems Design and Implementation*, pp. 469–482, 2017.
- [2] G. Pallis, “Cloud Computing: The New Frontier of Internet Computing,” *IEEE Internet Computing*, vol. 14, pp. 70–73, sep 2010.
- [3] P. M. Mell and T. Grance, “The NIST definition of cloud computing,” tech. rep., National Institute of Standards and Technology, Gaithersburg, MD, 2011.
- [4] Intricately, “2019 Intricately Cloud Market Share Report,” tech. rep., Intricately, 2019.
- [5] M. Armbrust, A. Fox, R. Griffith, A. D. Joseph, R. H. Katz, A. Konwinski, G. Lee, D. A. Patterson, A. Rabkin, I. Stoica, and M. Zaharia, “Above the Clouds: A Berkeley View of Cloud Computing,” *EECS Department, University of California, Berkeley*, pp. 1–25, feb 2009.
- [6] B. Power and J. Weinman, “Revenue growth is the primary benefit of the cloud,” *IEEE Cloud Computing*, vol. 5, pp. 89–94, jul 2018.
- [7] G. B. Berriman, E. Deelman, G. Juve, M. Rynge, and J. S. Vöckler, “The application of cloud computing to scientific workflows: A study of cost and performance,” *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, vol. 371, no. 1983, 2013.
- [8] A. Ruiz-Alvarez and M. Humphrey, “An automated approach to cloud storage service selection,” in *Proceedings of the 2nd international workshop on Scientific cloud computing - ScienceCloud '11*, (San Jose, California, USA), p. 39, ACM Press, 2011.

- [9] K. Bilal and A. Erbad, “Impact of multiple video representations in live streaming: A cost, bandwidth, and QoE analysis,” in *Proceedings - 2017 IEEE International Conference on Cloud Engineering, IC2E 2017*, pp. 88–94, 2017.
- [10] X. Pu, L. Liu, Y. Mei, S. Sivathanu, Y. Koh, and C. Pu, “Understanding performance interference of I/O workload in virtualized cloud environments,” in *Proceedings - 2010 IEEE 3rd International Conference on Cloud Computing, CLOUD 2010*, pp. 51–58, 2010.
- [11] A. Berl, E. Gelenbe, M. Di Girolamo, G. Giuliani, H. De Meer, M. Q. Dang, and K. Pentikousis, “Energy-Efficient Cloud Computing,” *The Computer Journal*, vol. 53, pp. 1045–1051, sep 2010.
- [12] L. Gkatzikis and I. Koutsopoulos, “Migrate or not? Exploiting dynamic task migration in mobile cloud computing systems,” *IEEE Wireless Communications*, vol. 20, no. 3, pp. 24–32, 2013.
- [13] N. J. Yadwadkar, B. Hariharan, J. E. Gonzalez, B. Smith, and R. H. Katz, “Selecting the best VM across multiple public clouds,” in *Proceedings of the 2017 Symposium on Cloud Computing - SoCC '17*, (Santa Clara, CA, USA), pp. 452–465, ACM Press, 2017.
- [14] P. Leitner and J. Cito, “Patterns in the Chaos - a Study of Performance Variation and Predictability in Public IaaS Clouds,” *ACM Transactions on Internet Technology*, vol. 16, pp. 1–23, apr 2014.
- [15] C. Davatz, C. Inzinger, J. Scheuner, and P. Leitner, “An Approach and Case Study of Cloud Instance Type Selection for Multi-Tier Web Applications,” in *Proceedings - 2017 17th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing, CCGRID 2017*, pp. 534–543, IEEE, may 2017.
- [16] C. Laaber, J. Scheuner, and P. Leitner, “Software microbenchmarking in the cloud. How bad is it really?,” *Empirical Software Engineering*, pp. 1–40, apr 2019.
- [17] J. Scheuner and P. Leitner, “A Cloud Benchmark Suite Combining Micro and Applications Benchmarks,” pp. 161–166, 2018.
- [18] J. Nocedal and S. Wright, *Numerical Optimization 2nd Ed.* No. 9781447122234 in Springer Series in Operations Research and Financial Engineering, Springer New York, 2006.
- [19] S. Venkataraman, Z. Yang, M. Franklin, B. Recht, and I. Nsdi, “Ernest : Efficient Performance Prediction for Large-Scale Advanced Analytics,” *NSDI'16 Proceedings of the 13th USENIX conference on Networked Systems Design and Implementation*, pp. 363–378, 2016.

- [20] J. Snoek, H. Larochelle, and R. P. Adams, “Practical Bayesian Optimization of Machine Learning Algorithms,” *Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 2*, pp. 2951–2959, 2012.
- [21] S. Agarwal, S. Kandula, N. Bruno, M.-C. Wu, I. Stoica, and J. Zhou, “Re-optimizing data-parallel computing,” in *NSDI’12 Proceedings of the 9th USENIX Symposium on Networked Systems Design and Implementation*, pp. 281–294, USENIX, 2012.

Appendices

Testing Summary

User Manual