

quiho: Automated Performance Regression Using Fine Granularity Resource Utilization Profiles

Ivo Jimenez

UC Santa Cruz

ivo.jimenez@ucsc.edu

Noah Watkins

UC Santa Cruz

nmwatkin@ucsc.edu

Michael Sevilla

UC Santa Cruz

msevilla@ucsc.edu

Jay Lofstead

Sandia National Laboratories

gflofst@sandia.gov

Carlos Maltzahn

UC Santa Cruz

carlosm@ucsc.edu

ABSTRACT

We introduce *quiho*, a framework used in automated performance regression tests. *quiho* discovers hardware and system software resource utilization patterns that influence the performance of an application. It achieves this by applying sensitivity analysis, in particular statistical regression analysis (SRA), using application-independent performance feature vectors to characterize the performance of machines. The result of the SRA, in particular feature importance, is used as a proxy to identify hardware and low-level system software behavior. The relative importance of these features serve as a performance profile of an application, which is used to automatically validate its performance behavior across revisions. We demonstrate that *quiho* can successfully identify performance regressions by showing its effectiveness in profiling application performance for synthetically induced regressions as well as several found in real-world applications.

CCS CONCEPTS

- Software and its engineering → Software performance; Software testing and debugging; Acceptance testing; Empirical software validation;
- Social and professional topics → Automation;

1 INTRODUCTION

Quality assurance (QA) is an essential activity in the software engineering process [1–3]. Part of the QA pipeline involves the execution of performance regression tests, where the performance of the application is measured and contrasted against past versions [4–6]. Examples of metrics used in regression testing are throughput, latency, or resource utilization over time. These metrics are compared and when significant differences are found, this constitutes a regression.

One of the main challenges in performance regression testing is defining the criteria to decide whether a change in an application’s performance behavior is significant, that is, whether a regression

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference’17, July 2017, Washington, DC, USA

© 2017 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

has occurred [7]. Simply comparing values (e.g. runtime) is not enough, even if this is done in statistical terms (e.g. mean runtime within a pre-defined variability range). Traditionally, this investigation is done by an analyst in charge of looking at changes, possibly investigating deeply into the issue and finally determining whether a regression exists.

When investigating a candidate of a regression, one important task is to find bottlenecks [8]. Understanding the effects in performance that distinct hardware and low-level system software¹ components have on applications is an essential part of performance engineering [9–11]. One common approach is to monitor an application’s performance in order to understand which parts of the system an application is hammering on [5]. Automated solutions have been proposed [7,12–14]. The general approach of these is to analyze logs and/or metrics obtained as part of the execution of an application in order to automatically determine whether a regression has occurred. Most of them do this by creating prediction models that are checked against runtime metrics. As with any prediction model, there is the risk of false/positive negatives.

In this work, we present *quiho* an approach aimed at complementing automated performance regression testing by using system resource utilization profiles associated to an application. A resource utilization profile is obtained using Statistical Regression Analysis² (SRA) where application-independent performance feature vectors are used to characterize the performance of machines. The performance of an application is then analyzed applying SRA to build a model for predicting its performance, using the performance vectors as the independent variables and the application performance metric as the dependant variable. The results of the SRA for an application, in particular feature importance, is used as a proxy to characterize hardware and low-level system utilization behavior. The relative importance of these features serve as a performance profile of an application, which is used to automatically validate its performance behavior across multiple revisions of its code base.

In this article, we demonstrate that *quiho* can successfully identify performance regressions. We show (Section 4) that *quiho* (1) obtains resource utilization profiles for application that reflect what their codes do and (2) effectively uses these profiles to identify induced regressions as well as other regressions found in real-world applications. The contributions of our work are:

¹Throughout this paper, we use “system” to refer to hardware, firmware and the operating system (OS).

²We use the term *Statistical Regression Analysis* (SRA) to differentiate between regression testing in software engineering and regression analysis in statistics.

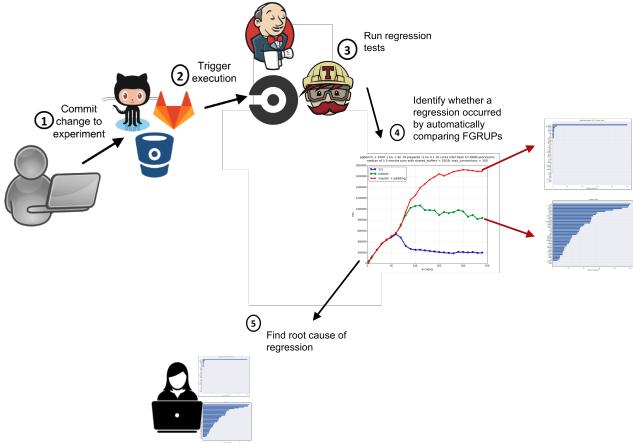


Figure 1: Automated regression testing pipeline integrating fine granularity resource utilization profiles (FGRUP). FGRUPs are obtained by *quiho* and can be used both, for identifying regressions, and to aid in the quest for finding the root cause of a regression.

- Insight: feature importance in SRA models (trained using these performance vectors) gives us a resource utilization profile of an application without having to look at the code.
- An automated end-to-end framework (based on the above finding), that aids analysts in identifying significant changes in resource utilization behavior of applications which can also aid in identifying root cause of regressions.
- Methodology for evaluating automated performance regression. We introduce a set of synthetic benchmarks aimed at evaluating automated regression testing without the need of real bug repositories. These benchmarks take as input parameters that determine their performance behavior, thus simulating different “versions” of an application.
- A negative result: ineffectiveness of resource utilization profiles for predicting performance using ensemble learning.

Next section (Section 2) shows the intuition behind *quiho* and how can be used to automate regression tests (Section 2). We then do a more in-depth description of *quiho* (Section 3), followed by our evaluation of this approach (Section 4). We briefly show how *quiho*'s resource utilization profiles can not be used to predict performance using some common machine learning techniques (Section 4.4). Section 5 reviews related work and we subsequently close with a brief discussion on challenges and opportunities enabled by *quiho* (Section 6).

2 MOTIVATION AND INTUITION BEHIND QUIHO

Fig. 1 shows the workflow of an automated regression testing pipeline and shows how *quiho* fits in this picture.

A regression is usually the result of observing a significant change in a performance metric of interest (e.g. runtime). At this point, an analyst will investigate further in order to find the root cause of the problem. One of these activities involves profiling an

application to see what's the pattern in terms of resource utilization. Traditionally, coarse-grained profiling (i.e. CPU-, memory- or IO-bound) can be obtained by monitoring an application's resource utilization over time. Fine granularity behavior allows application developers and performance engineers to quickly understand what they need to focus on while refactoring an application.

Obtaining fine granularity performance utilization behavior, for example, system subcomponents such as the OS memory mapping submodule or the CPU's cryptographic unit is usually time-consuming or requires implicates the use of more computing resources. This usually involves eyeballing source code, static code analysis, or analyzing hardware/OS performance counters.

An alternative is to infer fine granularity resource utilization behavior by comparing the performance of an application on platforms with different system performance characteristics. For example, if we know that machine A has higher memory bandwidth than machine B, and an application is memory-bound, then this application will perform better on machine A. There are several challenges with this approach:

1. We need to ensure that the software stack is the same on all machines where the application runs.
2. The amount of effort required to run applications on a multitude of platforms is not negligible.
3. It is difficult to obtain the performance characteristics of a machine by just looking at the hardware spec, so other more practical alternative is required..
4. Even if we could solve 3 and infer performance characteristics by just looking at the hardware specification of a machine, there is still the issue of not being able to correlate baseline performance with application behavior, since between two platforms is rarely the case where the change of performance is observed in only one subcomponent of the system (e.g. a newer machine doesn't have just faster memory sticks, but also better CPU, chipset, etc.).

The advent of cloud computing allows us to solve 1 using solutions like KVM [15] or software containers [16]. Chameleon-Cloud [17], CloudLab [18,19] and Grid5000 [20] are examples of bare-metal-as-a-service infrastructure available to researchers that can be used to automate regression testing pipelines for the purposes of investigating new approaches. These solutions to infrastructure automation coupled with DevOps practices [[21] ; htermann_devops_2012] allows us to address 2, i.e. to reduce the amount of work required to run tests.

Thus, the main challenge to inferring fine granularity resource utilization patterns (3 and 4) lies in quantifying the performance of the platform in a consistent way. One alternative is to look at the hardware specification and infer performance characteristics from this. As has been shown [needs-citation], this is not consistent. For example, the spec might specify that the machine has DDR4 memory sticks, with a theoretical peak throughput of 10 GB/s, but the actual memory bandwidth could be less (usually is, by a non-deterministic fraction of the advertised performance). *quiho* solves this problem by characterizing machine performance using microbenchmarks. These performance vectors are the “fingerprint” that characterizes the behavior of a machine [22].

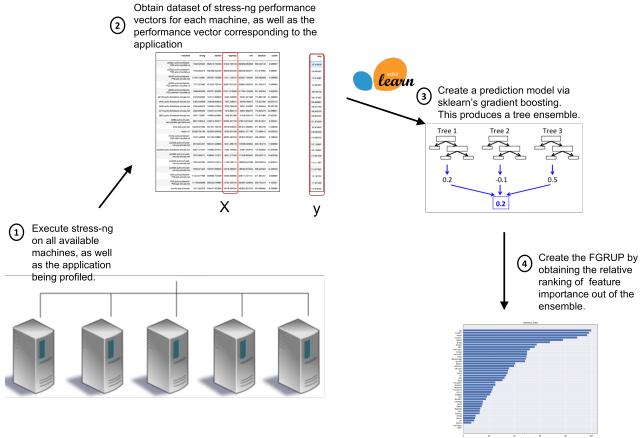


Figure 5: The workflow applied in order to obtain FGRUPs.

in computational simplicity of algorithms, presence of a closed-form solution, robustness with respect to heavy-tailed distributions, and theoretical assumptions needed to validate desirable statistical properties such as consistency and asymptotic efficiency. Some of the more common estimation techniques for linear regression are least-squares, maximum-likelihood estimation, among others.

scikit-learn [26] provides many of the previously mentioned techniques for building regression models. Another technique available in scikit-learn is gradient boosting [27]. Gradient boosting is a machine learning technique for regression and classification problems, which produces a prediction model in the form of an ensemble of weak prediction models, typically decision trees [28]. It builds the model in a stage-wise fashion like other boosting methods do, and it generalizes them by allowing optimization of an arbitrary differentiable loss function. This function is then optimized over a function space by iteratively choosing a function (weak hypothesis) that points in the negative gradient direction. Fig. 5 shows the process applied to obtaining FGRUPs for an application. We note that before creating the regression model, we normalize the data using the StandardScaler method from scikit-learn, which removes the mean from the dataset and scales the data to unit variance. Given that the bogo-ops-per-second metric does not quantify work consistently across stressors, we normalize the data in order to prevent some features from dominating in the process of creating the prediction models. In section Section 4 we evaluate the effectiveness of FGRUPs.

3.3 Using FGRUPs in Automated Regression Tests

As shown in Fig. 1 (step 4), when trying to determine whether a performance degradation occurred, FGRUPs can be used to compare differences between current and past versions of an application.

TODO: add algorithm

FGRUPs can also be used as a pointer to where to start with an investigation that looks for the root cause of the regression (Fig. 1, step 5). For example, if *memorymap* ends up being the most important feature, then we can start by looking at any code/libraries that make use of this subcomponent of the system. An analyst could

also trace an application using performance counters and look at corresponding performance counters to see which code paths make heavy use of the subcomponent in question.

4 EVALUATION

In this section we answer three main questions:

1. Can FGRUPs accurately capture application performance behavior? (Section 4.1)
2. Can FGRUPs work for identifying simulated regressions? (Section 4.2)
3. Can FGRUPs work for identifying regressions in real world software projects? (Section 4.3)
4. Can performance vectors be used to create performance prediction models? (Section 4.4)

Note on Replicability of Results: This paper adheres to The Popper Experimentation Protocol and convention⁶ [29], so experiments presented here are available in the repository for this article⁷. We note that rather than including all the results in the paper, we instead include representative ones for each section and leave the rest on the paper repository. Experiments can be examined in more detail, or even re-executed, by visiting the [source] link next to each figure. That link points to a Jupyter notebook that shows the analysis and source code for that graph. The parent folder of the notebook (following the Popper's file organization convention) contains all the artifacts and automation scripts for the experiments. All results presented here can be replicated⁸, as long as the reader has an account at Cloudbench (see repo for more details).

4.1 Effectiveness of FGRUPs to capture performance

In this subsection we show how FGRUPs can effectively describe the fine granularity resource utilization of an application with respect to a set of machines. Our methodology is:

1. Discover relevant performance features using the *quiho* framework.
2. Analyze source code to corroborate that discovered features are indeed the cause of performance differences.

We execute multiple applications for which fine granularity resource utilization characteristics we know in advance. These applications are redis [31], scikit-learn [26], and scca [32] and others. As a way to illustrate the variability originating from executing these applications on a heterogeneous set of machines, Fig. 6 shows boxplots of the four redis performance tests we execute.

In Fig. 7 we show four profiles side-by-side of four operations on redis, a popular open-source in-memory key-value database. These four tests are PUT, GET, LPOP and LPUSH. These benchmarks that test operations that put and get key-value pairs into the DB, and push/pop elements from a list stored in a key, respectively. The resource utilization profiles suggest that GET and PUT are memory intensive operations (first 3 stressors from each test, as shown in Tbl. 1). On the other hand, the profiles for LPOP and LPUSH look

⁶<http://falsifiable.us>

⁷<http://github.com/ivotron/quiho-popper>

⁸**Note to reviewers:** based on the terminology described in the ACM Badging Policy [30] this complies with the *Results Replicable* category. We plan to submit this work to the artifact review track too.

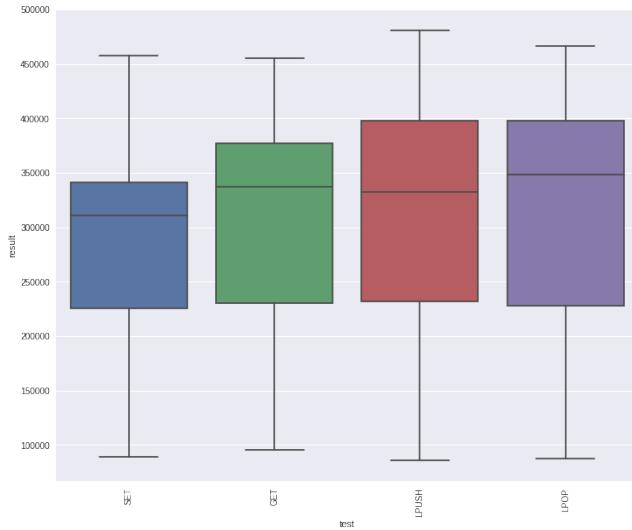


Figure 6: Variability. Y-axis is transactions per second.

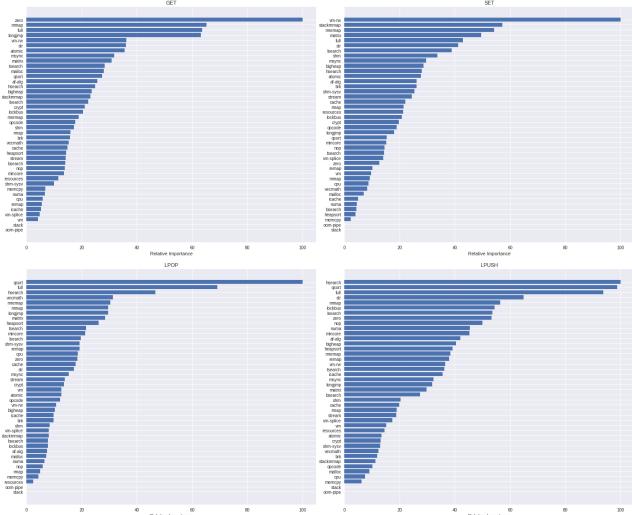


Figure 7: FGRUPs for four redis tests (PUT, GET, LPOP and LPUSH). These benchmarks that test operations that put and get key-value pairs into the DB, and push/pop elements from a list stored in a key, respectively.

different and they seem to have CPU intensive as the most important feature for this. If we look at the source code of redis, we can see why this is so. In the case of GET and PUT, these are memory intensive tasks. In the case of LPOP and LPUSH, these are routines that retrieve/replace the first element in the list, which is cpu-intensive and correlate with cpu-intensive stressors (such as `hsort` and `qsort`).

Fig. 8 shows the profile for one of the sklearn classification algorithm performance test. `scikit-learn` uses NumPy [33] internally, which is known to be memory-bound. `SSCA` on the other hand known to be CPU-bound.

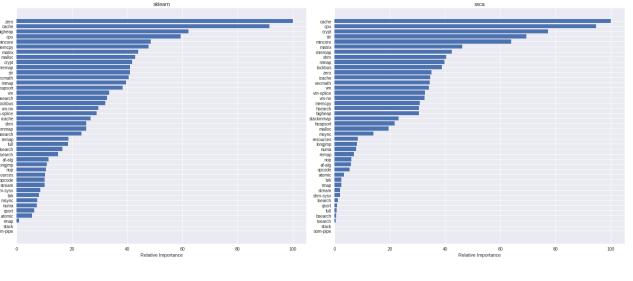


Figure 8: sklearn, sscfa.

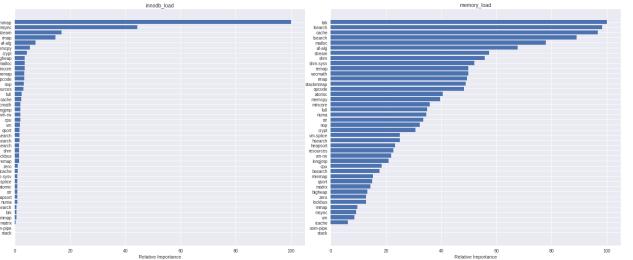


Figure 9: MariaDB with innodb and in-memory backends.

4.2 Simulating Regressions

In this section we test the effectiveness of `quiho` to detect performance simulations that are artificially induced. We induce regression by having a set of performance tests that take, as input, parameters that determine their performance behavior, thus simulating different “versions” of the same application. In total, we have 10 benchmarks for which we can induce several performance regressions, for a total of 30 performance regressions. For brevity, in this section we present results for two applications, MariaDB [34] and the STREAM-cycles.

The MariaDB test is based on the `mysqlslap` utility for stressing the database. In our case we run the load test, which populates a database whose schema is specified by the user. In our case, we have a fixed set of parameters that load a 10GB database. One of the exposed parameters is the one that selects the backend (storage engine in MySQL terminology). While the workload and test parameters are the same, the code paths are distinct and thus present different performance characteristics. The two engines we use in this case are `innodb` and `memory`. Fig. 9 shows the profiles of MariaDB performance for these two engines.

The next test is a modified version of the STREAM benchmark, which we refer to as STREAM-cycles. This version of STREAM introduces a `cycles` parameter that controls the number of times a STREAM operation is executed before reporting the time it took. In terms of the code, this adds an outer loop to each of the four different STREAM operations (`add`, `triad`, `copy`, `scale`), and loops as many times as the `cycles` parameter specifies. All STREAM tests are memory bound, so adding more cycles move the performance test from memory- to being cpu-bound; the higher the value of the `cycles` parameter, the more cpu-bound the test gets. Fig. 10 shows this behavior of all four tests across many machines.

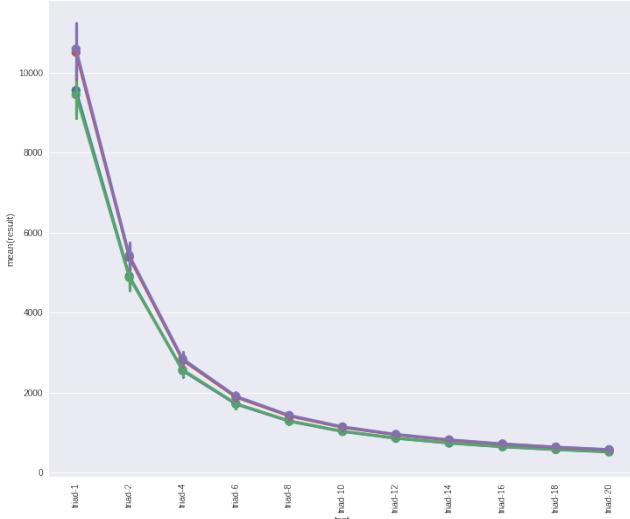


Figure 10: General behavior of the STREAM-cycles performance test. All STREAM tests are memory bound, so adding more cycles move the performance test from memory- to being cpu-bound; the higher the value of the cycles parameter, the more cpu-bound the test gets.

Fig. 11 shows the FGRUPs for the four tests. On the left, we see the “normal” resource utilization behavior of the “vanilla” version of STREAM (which corresponds to a value of 1 for the cycles parameter). As expected, the associated features (stressors) to these are from the memory/VM category. To the right, we see FGRUPs capturing the change in utilization behavior when cycles goes to its maximum value (20). In general FGRUPs do a good job of catching the simulated regression (which causes this application to be cpu-bound instead of memory-bound).

4.3 Real-world Scenario

In this section we show that *quiho* works with regressions that can be found in real software projects. It is documented that the changes made to the innodb storage engine in version 10.3.2 improves the performance in MariaDB, with respect to previous version 5.5.58. If we take the development timeline and invert it, we can treat 5.5.58 as if it was a “new” revision that introduces a performance regression. To show that this can be captured with FGRUPs, we use *mysqlslap* again and run the load test. Fig. 12 shows the corresponding FGRUPs. We can observe that the FGRUP generated by *quiho* can identify the difference in performance.

4.4 Using Performance Vectors to Predict Performance

As mentioned earlier, the set of performance vectors obtained as part of the generation of FGRUPs could be used to create prediction models that try to estimate the performance of an application using. Fig. 13 shows a plot with mean absolute percentage errors (MAPE) corresponding to the outcome of doing 1-cross-validation [35] across the distinct type of hardware architectures found in

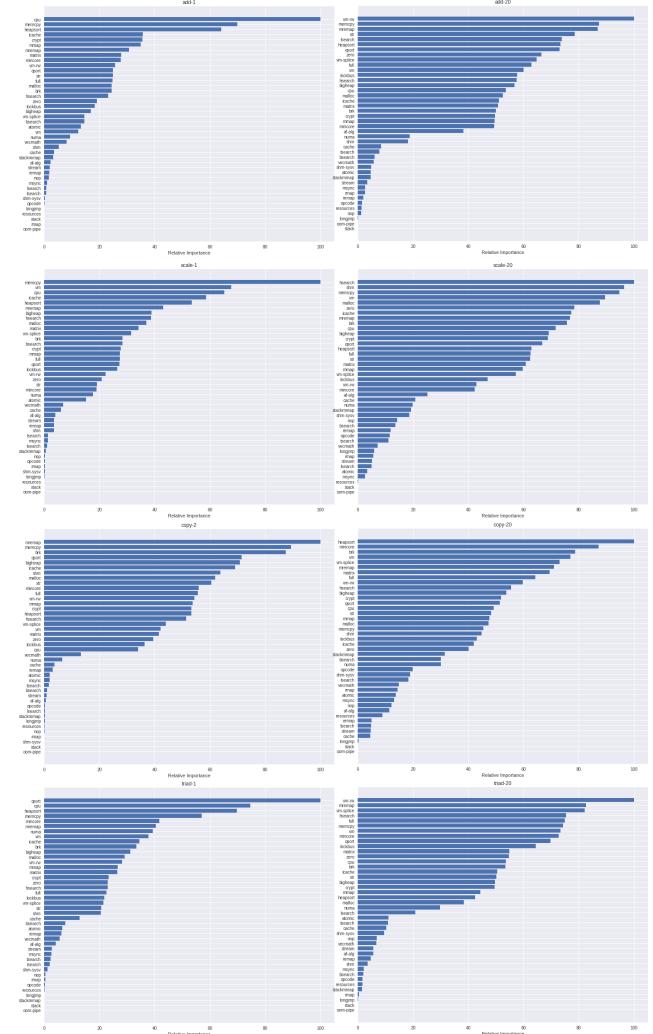


Figure 11: The FGRUPs for the four tests. We see that they capture the simulated regression (which causes this application to be cpu-bound instead of memory-bound).

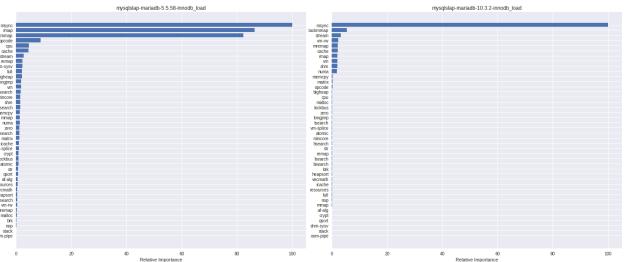


Figure 12: mariadb-10.0.3 vs. 5.5.

- Proceedings of the 6th ACM/SPEC International Conference on Performance Engineering*, 2015.
- [14] C. Heger, J. Happe, and R. Farahbod, "Automated Root Cause Isolation of Performance Regressions During Software Development," *Proceedings of the 4th ACM/SPEC International Conference on Performance Engineering*, 2013.
 - [15] A. Kivity, Y. Kamay, D. Laor, U. Lublin, and A. Liguori, "Kvm: The Linux virtual machine monitor," *Proceedings of the Linux symposium*, 2007.
 - [16] D. Merkel, "Docker: Lightweight Linux Containers for Consistent Development and Deployment," *Linux J.*, vol. 2014, Mar. 2014. Available at: <http://dl.acm.org/citation.cfm?id=2600239.2600241>.
 - [17] J. Mambretti, J. Chen, and F. Yeh, "Next Generation Clouds, the Chameleon Cloud Testbed, and Software Defined Networking (SDN)," *2015 International Conference on Cloud Computing Research and Innovation (ICCR)*, 2015.
 - [18] M. Hibler, R. Ricci, L. Stoller, J. Duerig, S. Guruprasad, T. Stack, K. Webb, and J. Lepreau, "Large-scale Virtualization in the Emulab Network Testbed," *USENIX 2008 Annual Technical Conference*, 2008. Available at: <http://dl.acm.org/citation.cfm?id=1404014.1404023>.
 - [19] R. Ricci and E. Eide, "Introducing CloudLab: Scientific Infrastructure for Advancing Cloud Architectures and Applications," *:login:*, vol. 39, 2014/December. Available at: <http://www.usenix.org/publications/login/dec14/ricci>.
 - [20] R. Bolze, F. Cappello, E. Caron, M. Daydé, F. Desprez, E. Jeannot, Y. Jégou, S. Lanteri, J. Leduc, N. Melab, G. Mornet, R. Namyst, P. Primet, B. Quetier, O. Richard, E.-G. Talbi, and I. Touche, "Grid'5000: A Large Scale And Highly Reconfigurable Experimental Grid Testbed," *Int J High Perform Comput Appl*, vol. 20, Nov. 2006.
 - [21] A. Wiggins, "The Twelve-Factor App" Available at: <http://12factor.net/>. Available at: <http://12factor.net/>.
 - [22] I. Jimenez, C. Maltzahn, J. Lofstead, A. Moody, K. Mohror, R. Arpac-Dusseau, and A. Arpac-Dusseau, "Characterizing and Reducing Cross-Platform Performance Variability Using OS-Level Virtualization," *2016 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, 2016.
 - [23] JW. Boys and D.R. Warn, "A Straightforward Model for Computer Performance Prediction," *ACM Comput Surv*, vol. 7, Jun. 1975.
 - [24] K. Kira and L.A. Rendell, "A Practical Approach to Feature Selection," *Proceedings of the Ninth International Workshop on Machine Learning*, 1992. Available at: <http://dl.acm.org/citation.cfm?id=645525.656966>.
 - [25] D.A. Freedman, *Statistical Models: Theory and Practice*, 2009.
 - [26] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and É. Duchesnay, "Scikit-learn: Machine Learning in Python," *J. Mach. Learn. Res.*, vol. 12, 2011. Available at: <http://www.jmlr.org/papers/v12/pedregosa11a.html>.
 - [27] P. Prettenhofer and G. Louppe, "Gradient Boosted Regression Trees in Scikit-Learn," Feb. 2014. Available at: <http://orbi.ulg.ac.be/handle/2268/163521>.
 - [28] J.H. Friedman, "Greedy Function Approximation: A Gradient Boosting Machine," *Ann. Stat.*, vol. 29, 2001.
 - [29] I. Jimenez, M. Sevilla, N. Watkins, C. Maltzahn, J. Lofstead, K. Mohror, A. Arpac-Dusseau, and R. Arpac-Dusseau, "The Popper Convention: Making Reproducible Systems Evaluation Practical," *2017 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, 2017.
 - [30] ACM, "Result and Artifact Review and Badging" Available at: <http://www.acm.org/publications/policies/artifact-review-badging>. Available at: <http://www.acm.org/publications/policies/artifact-review-badging>.
 - [31] J. Zawodny, "Redis: Lightweight key/value store that goes the extra mile," *Linux Mag.*, vol. 79, 2009.
 - [32] D.A. Bader and K. Madduri, "Design and Implementation of the HPCS Graph Analysis Benchmark on Symmetric Multiprocessors," *High Performance Computing - HiPC 2005*, 2005.
 - [33] S. van der Walt, S.C. Colbert, and G. Varoquaux, "The NumPy array: A structure for efficient numerical computation," *Comput. Sci. Eng.*, vol. 13, 2011.
 - [34] M. Widénius, "MariaDB SQL server project," *Ask Monty* Available at: <http://askmonty.org/wiki/index.php/MariaDB>. Available at: <http://askmonty.org/wiki/index.php/MariaDB>.
 - [35] R. Kohavi, "A study of cross-validation and bootstrap for accuracy estimation and model selection," *icai*, 1995.
 - [36] A. Liaw and M. Wiener, "Classification and regression by randomForest," *R News*, vol. 2, 2002.
 - [37] A. Crume, C. Maltzahn, L. Ward, T. Kroeger, and M. Curry, "Automatic generation of behavioral hard disk drive access time models," *2014 30th Symposium on Mass Storage Systems and Technologies (MSST)*, 2014.