

quiho: Automated Performance Regression Using Fine Granularity Resource Utilization Profiles

Ivo Jimenez

UC Santa Cruz

ivo.jimenez@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.

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 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.

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/nnnnnnnn.nnnnnnn>

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:

- 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.
- 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.

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.

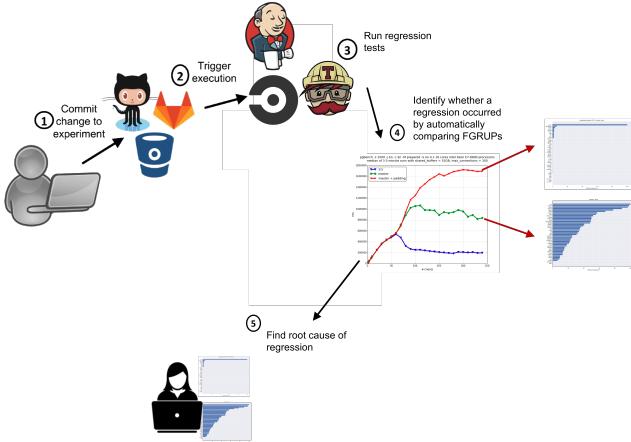


Figure 1: Automated regression testing pipeline integrating fine granularity resource utilization profiles (FGRUP). FGRUPs are obtained by *quiho* and can be used both, for identifying regressions, and to aid in the quest for finding the root cause of a regression.

- 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]. ChameleonCloud [17], CloudLab [18,19] and Grid5000 [20] are examples of bare-metal-as-a-service infrastructure available to researchers. DevOps [[21]; httermann_devops_2012] addresses 2). Thus, the main challenge with this approach 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 (Fig. 2). These performance vectors are the “fingerprint” that characterizes the behavior of a machine [22].

This performance vectors, obtained over a sufficiently large set of machines³, can serve as the foundation for building a prediction model of the performance of an application when executed on new (“unseen”) machines [23], a natural next step to take with a dataset like this. As we show in Section 4.4, this is not as good as we would expect.

However, building a prediction model has a utility. If we use these performance vectors to apply SRA and we focus on feature importance [24] of the created models, we can see that they give us fine granularity resource utilization patterns. In Fig. ??, we show the intuition behind why this is so. The performance of an application is determined by the performance of the subcomponents that get stressed the most by the application’s code. Thus, intuitively, if the performance of an application across multiple machines resembles

³As mentioned in Section 6, an open problem is to identify the minimal set of machines needed to obtaining meaningful results from SRA.

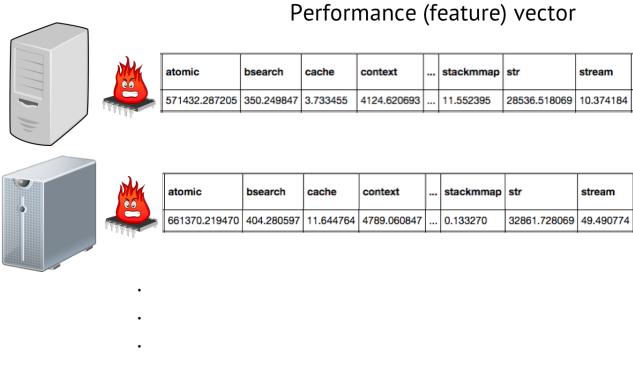


Figure 2: Performance vectors are obtained by executing a battery of microbenchmarks that quantify the performance of multiple subcomponents of a machine.

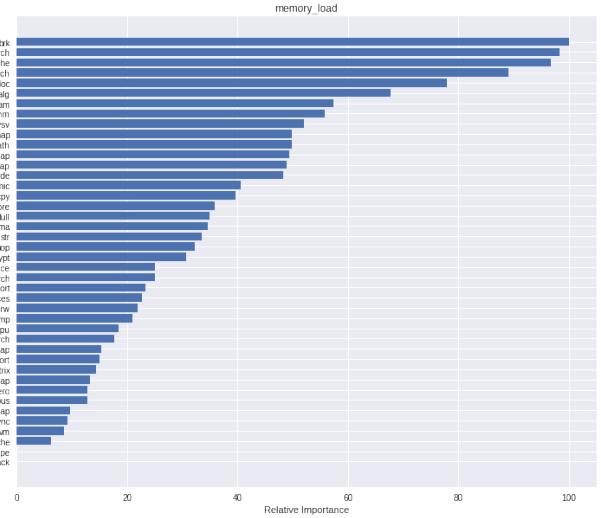


Figure 4: An example profile showing the relative importance of features for a particular application.

machine	af-alg	atomic	bigheap	brk	bsearch	cache	zlog
c220g1.quiho.schedock.net.ubuntu.genracks.net	7460.25556	680513.700493	37349.785748	69228.660098	369.462133	3.566657	97.018433
c220g2.quiho.schedock.net.ubuntu.genracks.net	7403.865310	68228.252459	36636.804958	682498.628441	370.318082	3.699681	106.484567
c220g2.quiho.schedock.net.ubuntu.genracks.net	11054.153981	68287.762689	11611.123370	303527.166342	308.383069	4.299998	72.519967
c220g2.quiho.schedock.net.ubuntu.genracks.net	7757.281868	451002.732744	26667.801230	506820.492253	251.355372	3.066660	94.091867
c220g2.quiho.schedock.net.ubuntu.genracks.net	7686.662426	450701.333302	27311.812889	511324.124985	251.933034	2.933282	89.742733
dc100g1.quiho.Schedock.emulab.net	4379.832862	164771.809846	2025.253988	32487.221988	173.985188	41.466890	188.191667
d430.quiho.Schedock.emulab.net	3454.033699	164628.906329	1907.548357	34439.766613	176.327592	38.000142	186.899667
d530.quiho.Schedock.emulab.net	4350.848073	16493.273432	2005.794649	30321.318062	175.592944	39.000188	186.231000
d710.quiho.Schedock.emulab.net	3456.095050	164447.914798	1915.892418	3397.885212	173.925279	40.566921	186.652333
d820.quiho.Schedock.emulab.net	4357.132681	16483.654692	1936.961068	31408.02412	173.915386	37.600261	186.823333
dc320g1.quiho.emulab.net.ubuntu.genracks.net	6611.338910	439114.302217	26323.321728	476418.274642	230.914851	3.333331	131.970000
dc320g2.quiho.emulab.net.ubuntu.genracks.net	13223.031616	657467.782752	43978.646205	691672.328268	371.884756	11.266639	62.879800
dc320g2.quiho.emulab.net.ubuntu.genracks.net	22839.765193	52303.564038	8025.801038	208221.311788	173.388412	19.533220	138.094000
dc320g2.quiho.emulab.net.ubuntu.genracks.net	7507.446968	53175.228881	48000.484393	767875.363427	299.498324	5.766648	118.982333
pc2400.quiho.emulab.net.ubuntu.genracks.net	9613.824407	380447.538905	9341.089105	175586.465862	236.164213	11.966682	107.150667
pc3000.quiho.Schedock.emulab.net	4381.731547	16490.357913	3422.164206	3791.053613	176.259883	37.933498	190.194667
pc3000.quiho.emulab.net.ubuntu.genracks.net	7670.092317	409655.727577	9621.317555	177405.855882	229.030712	14.800008	113.941333
pc3400.quiho.emulab.net.ubuntu.genracks.net	7802.481602	409704.016140	11132.209174	180905.497058	229.222454	14.800010	115.111667
pc3500.quiho.emulab.net.ubuntu.genracks.net	7850.021945	409757.668632	10076.285857	189780.573050	228.245323	14.238600	113.377667
r320.quiho.schedock.net.ubuntu.genracks.net	11054.775959	55098.754508	10022.069880	286115.721101	311.061257	3.966669	72.187767
r720.quiho.schedock.net.ubuntu.genracks.net	11148.859688	560523.049680	13725.433750	334997.529508	308.742572	4.100001	71.971000
scruffy.soe.ucsc.edu	7841.064578	605547.922968	19078.094834	382625.637476	250.386584	6.099968	112.679000

Variability patterns for a feature (across multiple machines), resembles the same variability pattern of application performance (across the same machines).

Figure 3: Intuition behind why feature importance implies resource utilization behavior. The variability patterns for a feature (across multiple machines), resembles the same variability pattern of application performance across the same machines. While this can be inferred by obtaining correlation coefficients, proper SRA is needed in order to create prediction models, as well as to obtain a relative rank of feature importances.

the performance of a microbenchmark, then we can say that the application is heavily influenced by that subcomponent.

If we rank features by their relative importance, we obtain what we call a fine granularity resource utilization profile (FGRUP), as shown in Fig. ??.

In the next section we show how these FGRUPs can be used in automated performance regression tests. Section 4 empirically validates this approach.

3 OUR APPROACH

In this section we do an in-depth description of *quiho*'s approach. We first describe how we obtain the performance vectors that characterize system performance. We then show how using these vectors we can feed SRA to build performance models for an application. Lastly, we describe how we obtain feature importance and how this represent a fine granularity resource utilization profile (FGRUP).

3.1 Performance Feature Vectors As System Performance Characterization

While the hardware and software specification can serve to describe the performance characteristics of a machine, the real performance characteristics can only feasibly⁴ be obtained by executing programs and capturing metrics. [can we show data for this? or add citation] The question then boils down to which programs should we use to characterize performance? Ideally, we would like to have many programs that execute every possible opcode mix so that we measure their performance. Since this is an impractical solution, an alternative is to create synthetic microbenchmarks that get as close as possible to exercising all the available features of a system.

*stress-ng*⁵ is a tool that is used to “stress test a computer system in various selectable ways. It was designed to exercise various physical subsystems of a computer as well as the various operating system kernel interfaces”. There are multiple stressors for CPU,

⁴One can get generate arbitrary performance characteristics by interposing a hardware emulation layer and deterministically associate performance characteristics to each instruction based on specific hardware specs. While possible, this is impractical (we are interested in characterizing “real” performance).

⁵<http://kernel.ubuntu.com/~cking/stress-ng>

machine	CPU	cores	run_cpus	threads	vcpus	board_name
\c228g1.quiho.schedock-pg8.wisc.cloudlab.us1*	\'Intel(R) Xeon(R) CPU E5-2630 v3 @ 2.40GHz\''	8	2	2	32	\'UCC-C228-M451\''
\c228g2.quiho.schedock-pg8.wisc.cloudlab.us1*	\'Intel(R) Xeon(R) CPU E5-2630 v3 @ 2.40GHz\''	8	2	2	32	\'UCC-C228-M451\''
\c6238.giho.schedock-pg8.apt.ululab.net1*	\'Intel(R) Xeon(R) CPU E5-2450 v8 @ 2.10GHz\''	8	1	2	16	\'PowerEdge R3281\''
\c6238.giho.schedock-pg8.clemon.cloudlab.us1*	\'Intel(R) Xeon(R) CPU E5-2660 v2 @ 2.20GHz\''	18	2	1	28	\'PowerEdge R3281\''
\c8228.giho.schedock-pg8.clemon.cloudlab.us1*	\'Intel(R) Xeon(R) CPU E5-2660 v2 @ 2.20GHz\''	18	2	1	28	\'PowerEdge R3281\''
\d1108.giho.schedock-emulab.net1*	\'Intel(R) Xeon(TM) CPU 3.0GHz\''	1	1	2	2	\'PowerEdge 2650\''
\d459.giho.schedock_emulab.net1*	\'Intel(R) Xeon(TM) CPU 3.0GHz\''	1	1	2	2	\'PowerEdge 2650\''
\d538.giho.schedock_emulab.net1*	\'Intel(R) Xeon(TM) CPU 3.0GHz\''	1	1	2	2	\'PowerEdge 2650\''
\d710.giho.schedock_emulab.net1*	\'Intel(R) Xeon(TM) CPU 3.0GHz\''	1	1	2	2	\'PowerEdge 2650\''
\d828.giho.schedock_emulab.net1*	\'Intel(R) Xeon(TM) CPU 3.0GHz\''	1	1	2	2	\'PowerEdge 2650\''
\d1368.giho.emulab-net.utahdc.genfracks.net1*	\'Intel(R) Xeon(R) CPU E5-2450 v8 @ 2.10GHz\''	8	2	2	32	\'ProLiant DL560e Gen8\''
\d1will.soe.ucsc.edu1*	\'Intel(R) Core(TM) i5-2400 CPU @ 3.10GHz\''	4	1	1	4	\'OptiPlex 790\''
\d5sdn-41.soe.ucsc.edu1*	\'Dual-Core AMD Opteron(tm) Processor 2212\''	2	2	1	4	\'AMD8321\''
\m18.giho.schedock-pg8.utah.cloudlab.us1*	\'Intel(R) Xeon(R) CPU D-1548 @ 2.80GHz\''	8	1	2	16	\'ProLiant #518 Server Cartridge\''
\pc2498.quiho.emulab-net.uky.emulab.net1*	\'Intel(R) Core(TM) i7 Quad CPU Q9550 @ 2.85GHz\''	4	1	1	4	\'Inspiron 5381\''
\pc3980.quiho.schedock_emulab.net1*	\'Intel(R) Xeon(TM) CPU 3.0GHz\''	1	1	2	2	\'PowerEdge 2650\''
\pc3986.quiho.emulab-net.uky.emulab.net1*	\'Intel(R) Core(TM) i7 Quad CPU Q6680 @ 2.40GHz\''	4	1	1	4	\'Inspiron 5381\''
\pc3948.quiho.emulab-net.uky.emulab.net1*	\'Intel(R) Core(TM) i7 Quad CPU Q6680 @ 2.40GHz\''	4	1	1	4	\'Inspiron 5381\''
\pc3950.quiho.emulab-net.uky.emulab.net1*	\'Intel(R) Core(TM) i7 Quad CPU Q6680 @ 2.40GHz\''	4	1	1	4	\'Inspiron 5381\''
\r128.giho.schedock-pg8.apt.emulab.net1*	\'Intel(R) Xeon(R) CPU E5-2450 v8 @ 2.10GHz\''	8	1	2	16	\'PowerEdge R3281\''
\r728.giho.schedock-pg8.apt.emulab.net1*	\'Intel(R) Xeon(R) CPU E5-2450 v8 @ 2.10GHz\''	8	1	2	16	\'PowerEdge R3281\''
\r728.soe.ucsc.edu1*	\'Intel(R) Xeon(R) CPU E5520 @ 2.40GHz\''	4	1	2	8	\'X870-QF1\''

Figure 5: List of different type of machines available on Cloudlab.

CPU cache, memory, OS, network and filesystem. Since we focus on system performance bandwidth, we execute the (as of version 0.07.29) 42 stressors for CPU, memory and virtual virtual memory stressors. A *stressor* is a function that loops a for a fixed amount of time (i.e. a microbenchmark), exercising a particular subcomponent of the system. At the end of its execution, *stress-ng* reports the rate of iterations executed for the specified period of time (referred to as *bogo-ops-per-second*).

Using this battery of stressors, we can obtain a performance profile of a machine (a performance vector). When this vector is compared against the one corresponding to another machine, we can quantify the difference in performance between the two at a per-stressor level. Every stressor (element in the vector) can be mapped to basic features of the underlying platform. For example, *stream* to memory bandwidth, *zero* to memory mapping, *qsort* to sorting data, and so on and so forth. However, the performance of a stressor in this set is *not* completely orthogonal to the rest. Fig. 6 shows a heat-map of Pearson correlation coefficients for performance vectors obtained by executing *stress-ng* on all the distinct machine configurations available in CloudLab [19] (Fig. 5 shows a summary of their hardware specs). As the figure shows, some stressors are slightly correlated (those near 0) while others show high correlation between them (in Section 4.4 we apply principal component analysis to this dataset).

3.2 System Resource Utilization Via Feature Importance in SRA

SRA is an approach for modeling the relationship between variables, usually corresponding to observed data points [25]. One or more independent variables are used to obtain a *regression function* that explains the values taken by a dependent variable. A common approach is to assume a *linear predictor function* and estimate the unknown parameters of the modeled relationships.

A large number of procedures have been developed for parameter estimation and inference in linear regression. These methods differ 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

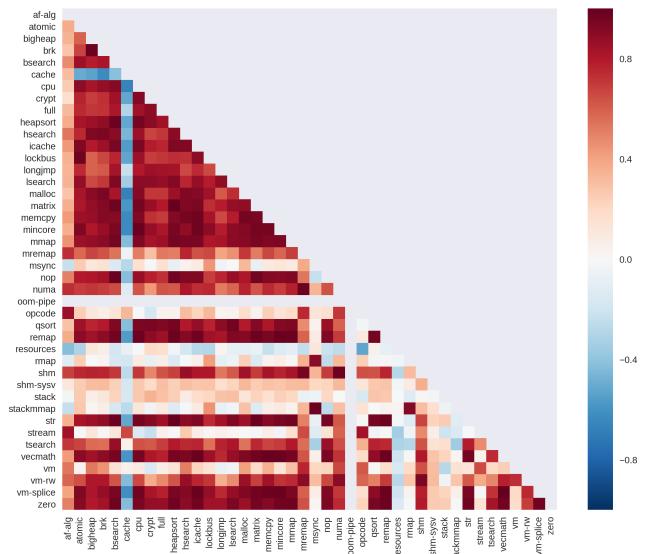


Figure 6: heat-map of Pearson correlation coefficients for performance vectors obtained by executing *stress-ng* on all the distinct machine configurations available in CloudLab.

the more common estimation techniques for linear regression are least-squares, maximum-likelihood estimation, among others.

sklearn [26] provides with many of the previously mentioned techniques for building regression models. Another technique available 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 function space by iteratively choosing a function (weak hypothesis) that points in the negative gradient direction. Fig. 7 shows the process applied to obtain FGRUPs for an application. In the next section we evaluate their effectiveness.

4 EVALUATION

In this section we answer:

- can FGRUPs accurately capture application performance behavior?
- can FGRUPs work for identifying simulated regressions?
- can FGRUPs work for identifying real-world regressions?

This paper adheres to The Popper Experimentation Protocol⁶ [29], so experiments presented here are available in the repository for this article⁷. 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, which points to an experiment and its artifacts.

⁶<http://falsifiable.us>

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

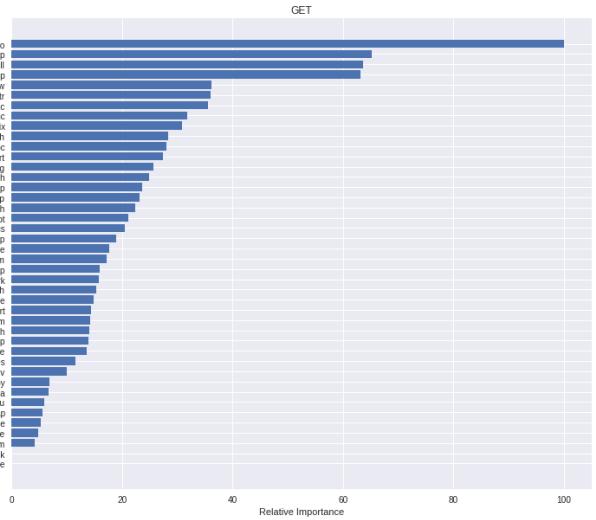
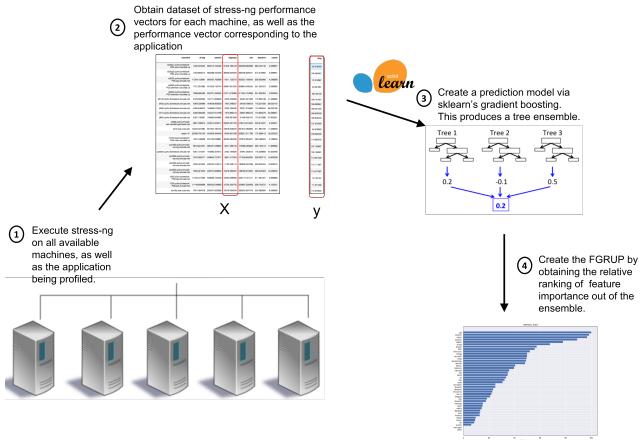


Figure 7: The workflow applied in order to obtain FGRUPs.

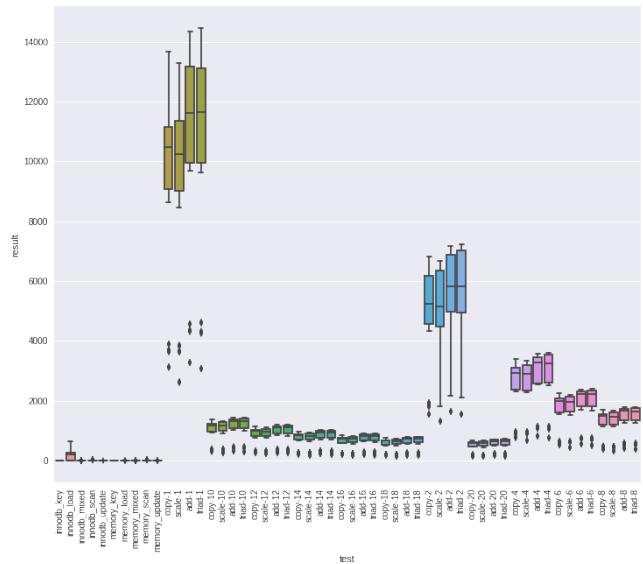


Figure 8: Variability.

4.1 Performance Profiles

We show they actually show the performance of known benchmarks.

Methodology - For every workload:

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

“Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu

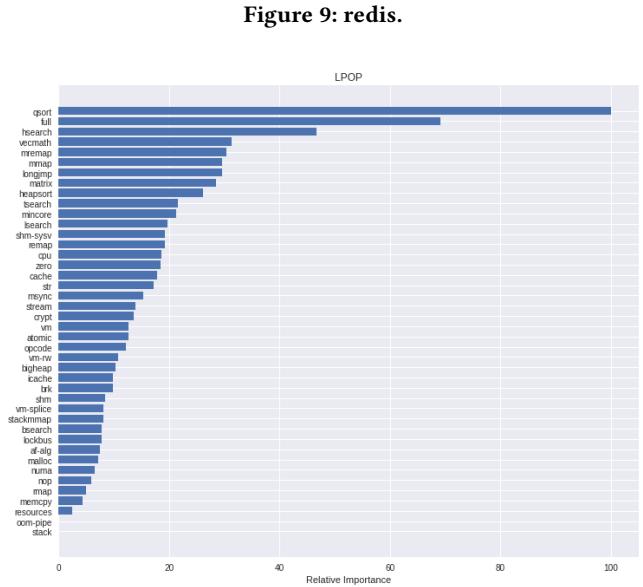


Figure 10: redis lpop.

fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.”

“Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.”

"Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in

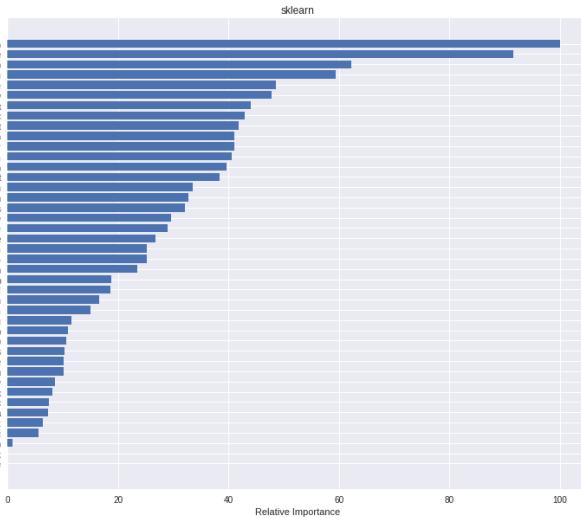


Figure 11: sklearn

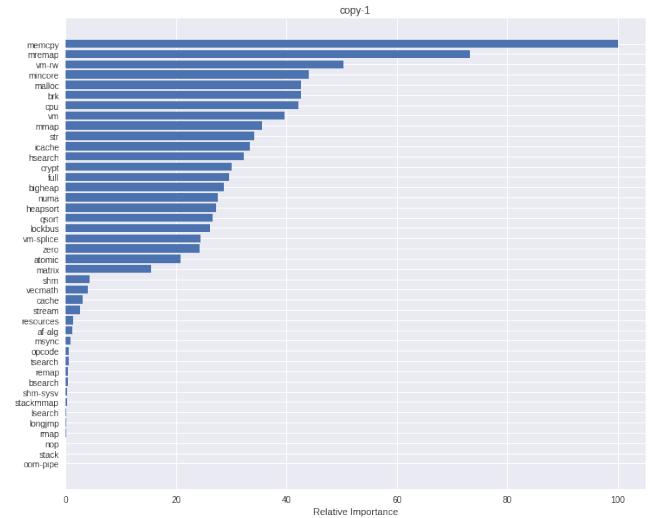


Figure 13: stream copy-1

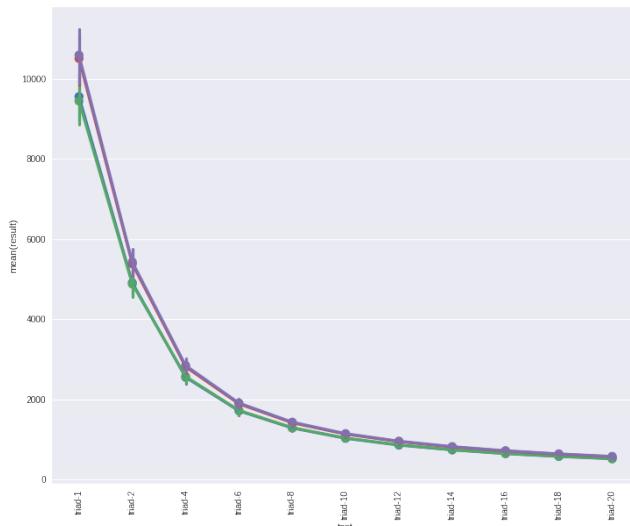


Figure 12: stream cycles

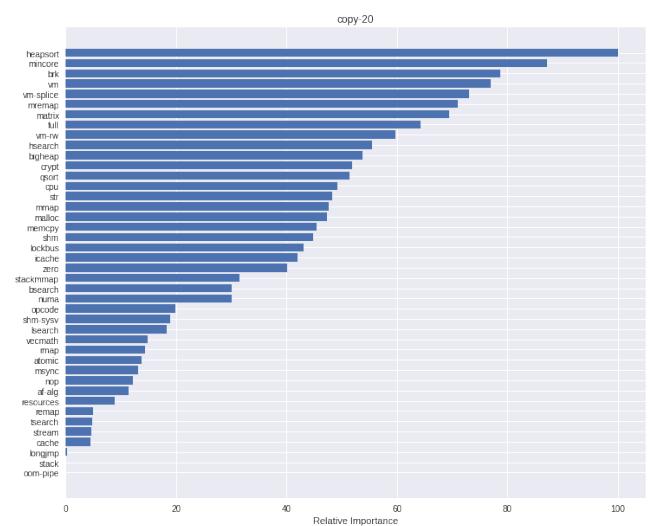


Figure 14: stream copy-20

reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.”

4.2 Simulating Regressions

We show that if we simulate regressions, then *quiho* identifies them correctly.

“Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.”

“Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.”

text in between “Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.”

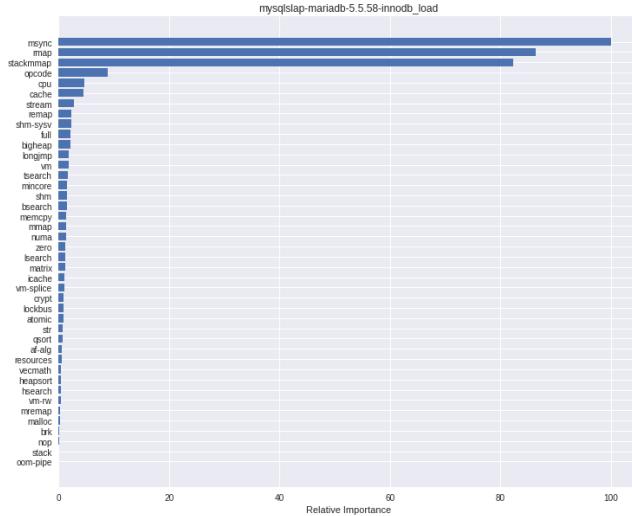


Figure 15: mariadb-5.5.58.

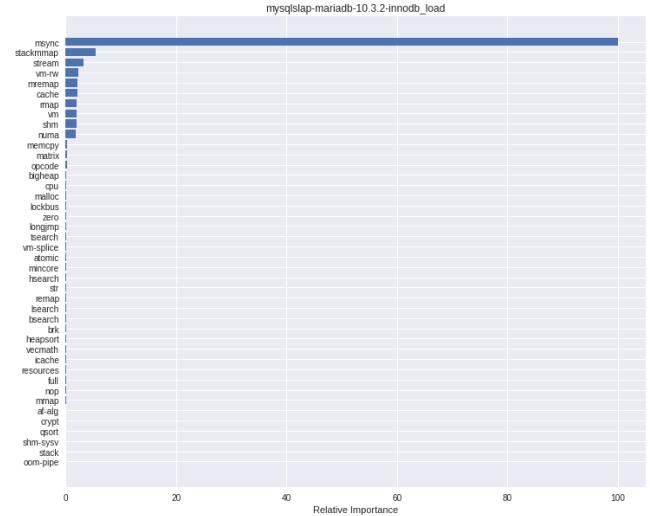


Figure 16: mariadb-10.0.3

mariadb innodb.

"Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum."

mariadb memory.

Due to lack of space we can not include the remaining benchmarks. In total, we have 10 benchmarks for which we can induce several performance regressions, for a total of 30 performance regressions.

4.3 Real-world Scenario

We show that *quiho* works with a real software project.

Use of FGRUPs: FGRUPs aid in performance engineering. When analyzing any performance degradation of an application, then the feature importance can be used as a pointer to where to start with the investigation. For example, if *memorymap* ends up being the most important feature, then we can start by looking at any code/libraries related to this functionality, or looking at corresponding performance counters (if available). In this case, we would look for the code.

Due to lack of space, we cannot present results for 5 other applications where we detect regressions successfully. These are ZLog, Apache Commons Math, GCC, PostgreSQL, Redis, Apache Web Server and MongoDB.

4.4 *quiho* cannot predict performance

We show how *quiho* does not do a good job at predicting performance.

text in between
text in between

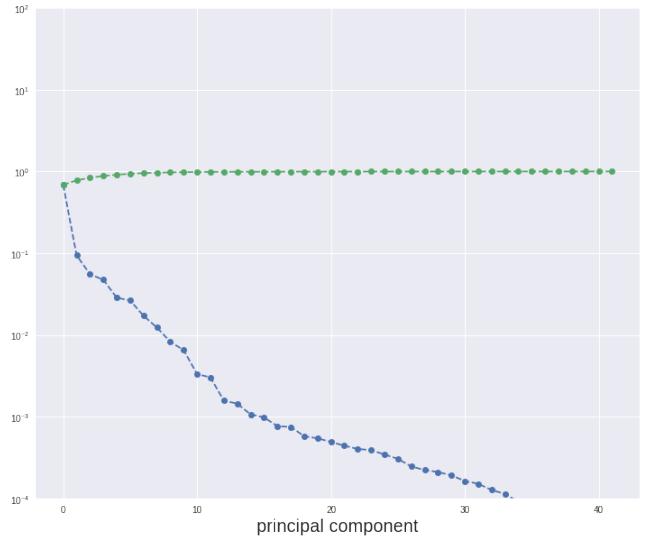


Figure 17: Variability reduction per subcomponent in PCA.

5 RELATED WORK

5.1 Anomaly Detection and Bottleneck Identification

It's been used in bottleneck detection [8]. TODO: mention briefly how it is used.

5.2 Automated Regression Testing

In [13], they use it to detect regressions using a dataset of performance counters.

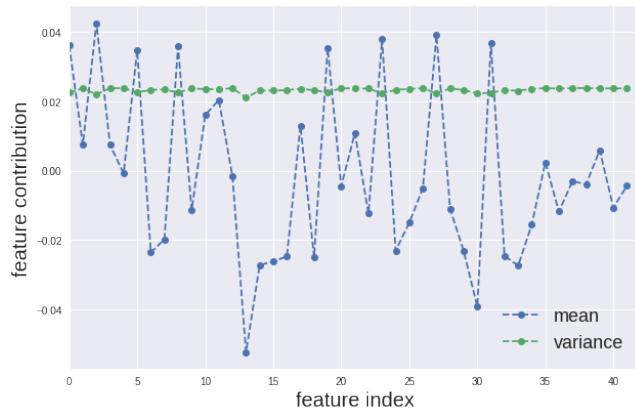


Figure 18: Variability reduction per index.

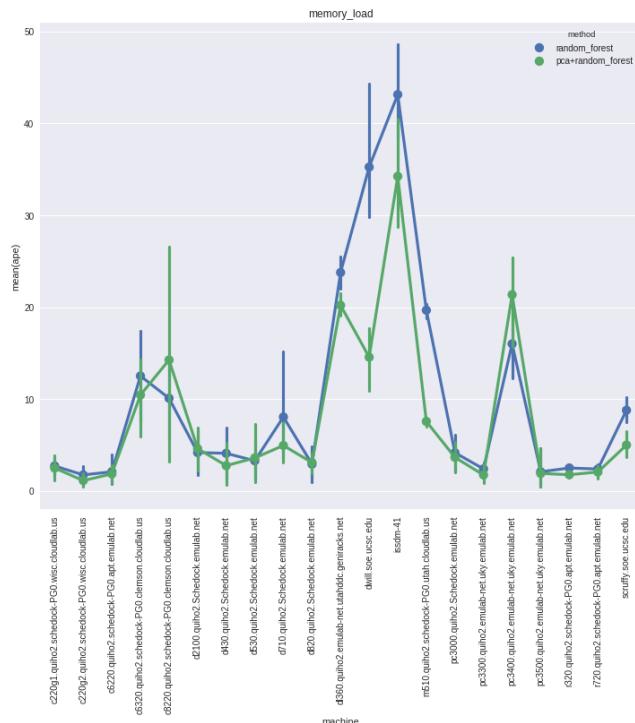


Figure 19: Mean Absolute Percentage Error of cross-validation.

6 CONCLUSION AND FUTURE WORK

- Main draw-back of this technique is that we need to run on multiple machines.
 - We used stress-ing but this is not the only thing we can use. Ideally, we would extend this battery of tests so that we have more “coverage” of the distinct subcomponents of a system.

In the not-so-distant future:

- multi-node

- minimum number of machines?
 - single machine?
 - long-running (multi-stage) applications. e.g. a web-service or big-data application with multiple stages. In this case, we would define windows of time and we would apply quiho to each. The challenge: how do we automatically get the windows rightly placed.

Acknowledgments: This work was partially funded by the Center for Research in Open Source Software⁸, Sandia National Laboratories and NSF Awards #1450488.

REFERENCES

- [1] G.J. Myers, C. Sandler, and T. Badgett, *The Art of Software Testing*, 2011.
 - [2] A. Bertolini, "Software Testing Research: Achievements, Challenges, Dreams," *2007 Future of Software Engineering*, 2007.
 - [3] B. Beizer, *Software Testing Techniques*, 1990.
 - [4] J. Dean and L.A. Barroso, "The tail at scale," *Commun ACM*, vol. 56, Feb. 2013.
 - [5] B. Gregg, *Systems Performance: Enterprise and the Cloud*, 2013.
 - [6] F.I. Vokolos and E.J. Weyuker, "Performance Testing of Software Systems," *Proceedings of the 1st International Workshop on Software and Performance*, 1998.
 - [7] L. Cherkasová, K. Ozonat, N. Mi, J. Symons, and E. Smirni, "Anomaly? Application change? Or workload change? Towards automated detection of application performance anomaly and change," *2008 IEEE International Conference on Dependable Systems and Networks With FTCS and DCC (DSN)*, 2008.
 - [8] O. Ibdumoye, F. Hernández-Rodríguez, and E. Elmroth, "Performance Anomaly Detection and Bottleneck Identification," *ACM Comput Surv*, vol. 48, Jul. 2015.
 - [9] G. Jin, L. Song, X. Shi, J. Scherpelz, and S. Lu, "Understanding and Detecting Real-world Performance Bugs," *Proceedings of the 33rd ACM SIGPLAN Conference on Programming Language Design and Implementation*, 2012.
 - [10] S. Han, Y. Dang, S. Ge, D. Zhang, and T. Xie, "Performance Debugging in the Large via Mining Millions of Stack Traces," *Proceedings of the 34th International Conference on Software Engineering*, 2012.
 - [11] M. Jovic, A. Adamoli, and M. Hauswirth, "Catch Me if You Can: Performance Bug Detection in the Wild," *Proceedings of the 2011 ACM International Conference on Object Oriented Programming Systems Languages and Applications*, 2011.
 - [12] Z.M. Jiang, "Automated Analysis of Load Testing Results," *Proceedings of the 19th International Symposium on Software Testing and Analysis*, 2010.
 - [13] W. Shang, A.E. Hassan, M. Nasser, and P. Flora, "Automated Detection of Performance Regressions Using Regression Models on Clustered Performance Counters," *Proceedings of the 6th ACM/SPEC International Conference on Performance Engineering*, 2015.
 - [14] C. Heger, J. Happé, 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. Liguri, "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.
 - [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 (ICCCR)*, 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.
 - [19] R. Ricci and E. Eide, "Introducing CloudLab: Scientific Infrastructure for Advancing Cloud Architectures and Applications," *:login*, vol. 39, 2014/December.
 - [20] R. Bolze, F. Cappello, E. Caron, M. Daydé, F. Despres, 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"
 - [22] I. Jimenez, C. Maltzahn, J. Lofstead, A. Moody, K. Mohror, R. Arpacı-Dusseau, and A. Arpacı-Dusseau, "Characterizing and Reducing Cross-Platform Performance

⁸<http://cross.ucsc.edu>

- Variability Using OS-Level Virtualization," *2016 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, 2016.
- [23] J.W. Boysen 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.
- [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.
- [27] P. Prettenhofer and G. Louppe, "Gradient Boosted Regression Trees in Scikit-Learn," Feb. 2014.
- [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. Arpacı-Dusseau, and R. Arpacı-Dusseau, "The Popper Convention: Making Reproducible Systems Evaluation Practical," *2017 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, 2017.