

# Analyzing the Impact of Workloads on Modeling the Performance of Configurable Software Systems

Stefan Mühlbauer  
Leipzig University

Florian Sattler  
Saarland University  
Saarland Informatics Campus

Christian Kaltenecker  
Saarland University  
Saarland Informatics Campus

Johannes Dorn  
Leipzig University

Sven Apel  
Saarland University  
Saarland Informatics Campus

Norbert Siegmund  
Leipzig University  
ScaDS.AI Dresden/Leipzig

**Abstract**—Modern software systems often exhibit numerous configuration options to tailor them to user requirements, including the system’s performance behavior. Performance models derived via machine learning are an established approach for estimating and optimizing configuration-dependent software performance. Most existing approaches in this area rely on software performance measurements conducted with a single workload (i.e., input fed to a system). This single workload, however, is often not representative of a software system’s real-world application scenarios. Understanding to what extent configuration and workload—individually and combined—cause a software system’s performance to vary is key to understand whether performance models are generalizable across different configurations and workloads. Yet, so far, this aspect has not been systematically studied.

To fill this gap, we conducted a systematic empirical study across 25 258 configurations from nine real-world configurable software systems to investigate the effects of workload variation at system-level performance and for individual configuration options. We explore driving causes for workload–configuration interactions by enriching performance observations with option-specific code coverage information.

Our results demonstrate that workloads can induce substantial performance variation and interact with configuration options, often in *non-monotonous* ways. This limits not only the generalizability of single-workload models, but also challenges assumptions for existing transfer-learning techniques. As a result, workloads should be considered when building performance prediction models to maintain and improve representativeness and reliability.

## I. INTRODUCTION

Most modern software systems can be customized by means of configuration options enabling desired functionality or tweaking non-functional aspects, such as performance or energy consumption. The relationship between configuration choices and their influence on performance has been extensively studied in the literature [1]–[10]. The backbone of performance estimation are prediction models that map a given configuration to the estimated performance value. Learning performance models relies on a training set of configuration-specific performance measurements. In state-of-the-art approaches, observations usually rely on only single-workload measurements that aim at reflecting performance behavior of a typical real-world application scenario.

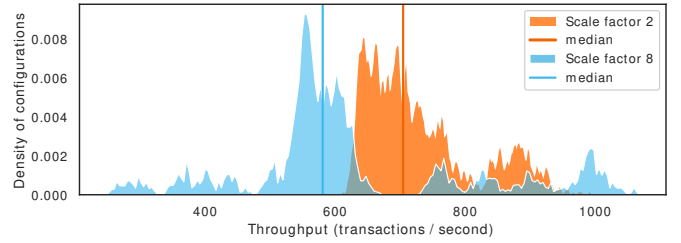


Figure 1: Throughput distribution of 1954 configurations of the database system H2 for the TPC-C benchmark at two different scale factors.

It is almost folklore that choice of the workload (i.e., the input fed to the software system) influences the performance of software systems in different ways [11], as has been shown for the domains of SAT solving [12], [13], compilation [14], [15], video transcoding [16], [17], data compression [18], and code verification [19]. Beside apparent interactions, such as performance scaling with the size of a workload, qualitative aspects can result in intricate and inadvertent performance variations.

Take as an example two performance throughput distributions across the configuration space of the database system H2 in Figure 1. Here, the exact same configurations on two different parameterizations of the benchmark TPC-C have been measured. In this setting, the scale factor controls the modeled number of warehouses.

While for most configurations, throughput decreases when the scale factor is increased some configurations achieve even a higher throughput<sup>1</sup>. This example illustrates that configuration-dependent performance can be highly sensitive to workload variation and that the performance behavior under different workloads can change in unexpected ways. In turn, this can render performance models based on a single workload useless, unless the configuration options’ sensitivity to workloads is accounted for.

<sup>1</sup>A similar workload-specific performance distribution was described by Pereira et al. for the video encoder x264 [17].

To address this limitation, two different approaches have been pursued in the literature: performance modeling (1) based on existing knowledge [20]–[24], and (2) for a combined configuration-workload problem space [19], [25].

The first approach relies on transfer-learning techniques, in which, given an existing performance model, in a second step only the differences to a new environment or workload are learned. A transfer function encodes which configuration options’ influences on performance are sensitive to workload variation. While transfer learning is an effective strategy that is not limited to varying workloads [20]–[24], its main limitation is that the transfer function is specific to the differences between two environments.

In contrast to transfer learning, a second and more generalist approach is to consider the input fed to a software system as a further dimension for modeling performance. A workload is characterized by properties that—individually or in conjunction with software configuration options—influence performance. For such a strategy to work, one requires in-depth knowledge of the characteristics of a workload that influence performance, let alone these characteristics can be mathematically modeled at all. This strategy has been effectively tested for a variety of application domains, such as program verification [19] and high-performance computing [25]. However, the added complexity comes at substantial cost. Not only does it require substantially more measurements, one often lacks knowledge of which performance-relevant characteristics best describe a workload (e.g., what makes a program hard to verify or optimize).

The existing body of research [1]–[10], [26]–[29] confirms the prevalence and importance of the influence of the workload on the performance of software systems. All these works are aware of the workload dimension as a factor of performance variation, yet little is known about the quality and driving factors of the *interplay* between configuration options and workloads. Are workloads and configurations as two factors influencing software performance orthogonal and can be treated independently, or does their interplay give rise to intricate and inadvertent performance behavior? For example, varying the workload had a non-uniform effect on different configurations of the database system H2 (cf. Fig. 1), suggesting a specific interaction between certain configuration options and the workload.

We have conducted an empirical study that sheds light on whether and how choices of configuration and workload interact with regard to performance. Specifically, we have analyzed 25 258 configurations from nine configurable real-world software systems to obtain a broad picture of the interaction of configuration and workload when learning performance models and estimating performance (i.e., response time). Aside from studying the sole effects of workload variation on performance behavior, we explore *what* drives the interaction between workload and configuration. To this end, we enrich performance observations with corresponding coverage data to understand workload variation with respect to the executed code.

We found that varying the workload can influence configuration-dependent software performance in different ways, including non-linear and non-monotonous effects. As a key take-away, we provide empirical evidence that single-workload approaches do not generalize across workload variations and that even existing transfer-learning techniques are too limited to address non-monotonous performance variations induced by qualitative workload changes. We demonstrate how coverage testing can outline a path to screen for workload-sensitive configuration options.

To summarize, we make the following contributions:

- An empirical study of 25 258 configurations from nine configurable software systems on whether and how interactions of workload and configuration affect software performance;
- A detailed analysis illustrating that (1) system-level performance, and (2) the performance influence of individual configuration options can be sensitive to workload variation and often exhibit a *non-monotonous* relationship, caused by variation in the execution of option-specific code;
- A critical reflection of the suitability of single-workload models for predicting configuration-dependent performance and assumptions of state-of-the-art transfer-learning approaches in this area;
- An archived repository on [zenodo.org](https://zenodo.org)<sup>2</sup> with supplementary material, including performance and coverage measurements, configurations, and an interactive dashboard for data exploration to *reproduce* all analyses and additional visualizations left out due to space limitations.

## II. PRELIMINARIES AND RELATED WORK

Software performance emerges from a variety of factors including configuration, workload, and hardware setup. In what follows, we revisit work that models the relationship between such factors (individually or in combinations) and software performance.

### A. Performance Prediction Models

Configurable software systems is an umbrella term for any kind of software system that exhibits configuration options to customize its functionality [30]. While the primary purpose of configuration options is to select and tune functionality, each configuration choice may also have implications on non-functional properties (e.g., execution time or memory usage)—be it intentional or not. Performance prediction models approximate non-functional properties, such as execution time or memory usage, as a function of software configurations  $c \in C$ , formally,  $\Pi : C \rightarrow \mathbb{R}$ .

Such *black-box* models do not rely on an understanding of the internals of a configurable software system, but are learned from a training set of configuration-specific performance observations. In this vein, finding configurations with optimal performance [26]–[29] and estimating the performance

<sup>2</sup><https://doi.org/10.5281/zenodo.7504284>

for arbitrary configurations of the configuration space is an established line of research [1]–[10]. Over the decade, several different learning and modeling techniques have shown to be effective to learn configuration-dependent software performance, including probabilistic programming [1], multiple linear regression [2], classification and regression trees [5]–[7], Fourier learning [8], [9], and deep neural networks [3], [4], [10]. The set of configurations for training can be sampled from the configuration space using a variety of different sampling techniques [31], [32]. All sampling techniques aim at yielding a representative sample, either by covering the main effects of configuration options and interactions among them [33] or by sampling uniformly from the configuration space [29], [34]. Most sampling techniques share the perspective of treating a configurable software system as a black-box model at application-level granularity. Recent work has incorporated feature location techniques to guide sampling effort towards relevant configuration options [35], [36] or to model non-functional properties at finer granularity [37], [38].

### B. Varying Workloads

When assessing the performance of a software system, one asks how well a certain *operation* is executed, or, phrased differently, how well an *input fed to the software system* is processed. In the context of this study, we will refer to such inputs as *workloads*. By nature, the workload of a software system is application-specific, such as a series of queries and transactions fed to a database system or a sequence of raw image files processed by a video encoder. Workloads can often be distinguished by the characteristics they exhibit, such as the amount and type of data to be processed (text, binary data).

In practice, a useful workload for assessing performance should closely resemble the real-world scenario that the system under test will be deployed in. To achieve this, a well-defined and widely employed technique in performance engineering is workload characterization [39], [40]. To select a representative workload, it is imperative to explore workload characteristics and validate a workload with real-world observations. This can be achieved by constructing workloads from usage patterns [41] or by increasing the workload coverage using a mix of different workloads rather than a single one [42].

While workload characterization and benchmark construction is domain-specific, there are numerous examples of this task being driven by community efforts. For instance, the non-profit organizations Transaction Processing Performance Council (TPC) and Standard Performance Evaluation Corporation (SPEC) provide large bodies of benchmarks for data-centric applications and across different domains, respectively.

### C. Workloads and Performance Prediction

Different approaches have been proposed to tackle the problem of workload sensitivity in performance prediction.

*a) Workload-aware Performance Modeling:* Extending on workload characterization (cf. Section II-B), a strategy that embraces workload diversity is to incorporate workload characteristics into the problem space of a performance prediction

model. Here, performance is modeled as a function of both the configuration options *explicitly* exhibited by the software system as well as the workload characteristics, formally  $\Pi : C \times W \rightarrow \mathbb{R}$ . The combined problem space enables learning performance models that generalize to workloads that exhibit characteristics denoted by  $W$  since we can screen for performance-relevant combinations of options and workload characteristics. This domain-specific strategy has been successfully applied to domains such as program verification [19], algorithm selection [43], or the parametrization of the Java microbenchmark harness [44]. In these instances, the characteristics (varying aspects of a workload) are explicitly specified and do not require further characterization.

Its main disadvantages are twofold: The combined problem space (configuration and workload dimension) requires substantially more observations to screen for identifying performance-relevant options, characteristics, and interactions thereof. In addition, previous work found that only few configuration options are sensitive to workload variation [21]. That is, the problem of identifying meaningful, but sparse predictors is exacerbated since one must not only identify performance-relevant configuration options but also workload-sensitive ones. It is not possible to find such a characterization in every case. Even worse, a chosen characterization can be wrong and omit important factors or overestimate unimportant factors.

At large, the notion of the influence of workloads on configuration-dependent performance remains the exception in the literature: While a study related to ours explores and confirms the presence of interactions between the workload and configuration options [45], only few researchers even consider this dimension of the problem space.

*b) Transfer Learning for Performance Models:* Another strategy for workload-aware performance prediction builds on the fact that, across different workloads, only few configuration options are in fact workload sensitive [21]. One first trains a model on a standard workload and, subsequently, adapts it to a different workload of choice. Contrary to a generalizable workload-aware model, transfer-learning strategies focus on approximating a transfer function that, without characterizing the workload, encodes the information on which configuration options are sensitive to differences between a source and target pair of workloads. Training a workload-specific model and adapting it on demand provides an effective means to reuse performance models, which is not only limited to workloads [20], [23], [24], [46]. The main shortcoming of transfer learning is that it does not generalize to arbitrary workloads, since a transfer function is tailored to a specific target workload. Basically, one trades generalizability for measurement cost, because learning a transfer function requires substantially fewer training samples.

While both directions (workload-aware performance modeling and transfer learning) are effective means to handle workload sensitivity, to the best of our knowledge, there is no *systematic* assessment of the factors that drive the interaction between configuration and workload with regard to performance. Understanding scenarios that are associated

with or even cause incongruent performance influences across workloads (1) help practitioners to employ established analysis techniques more effectively and (2) motivate researchers to devise analysis techniques dedicated to such scenarios.

### III. STUDY DESIGN

In what follows, we describe our research questions and measurement setup. We make all performance measurement data, configurations, workloads, and learned performance models available on the paper’s companion Web site.

#### A. Research Questions

The first two research questions are concerned with the workload sensitivity of the studied software systems’ performance behavior. We first take a look at the entire system (RQ<sub>1</sub>) and its configurations and, subsequently, to individual configuration options (RQ<sub>2</sub>). In Sec. V, we explore possible driving factors and indicators for workload-specific performance variation of configuration options (RQ<sub>3</sub>).

1) *Performance Variation Across Workloads*: Performance variation may arise from workload variation [11]. In a practical setting, the question arises whether, and if so, to what extent an existing workload-specific performance model is representative of the performance behavior of also other workloads. That is, can a model estimating the performance of different configurations be reused for the same software system but run with a different workload? Clearly, it depends. But, analyzing systematically how the degree of similarity of workloads and corresponding performance behaviors varies across the configuration space provides insights into the extent the strategies of transferring performance models (outlined in Section II-C) might be applicable. To this end, we formulate the following research question:

RQ<sub>1</sub> *To what extent does performance behavior vary across workloads and configurations?*

2) *Option Influence Across Workloads*: The global performance behavior emerges from the influences of several individual options and their interaction as well as the combined influence with the workload on performance. To understand which configuration options are driving performance variation, in general, and which are workload sensitive, in particular, we state the following research question:

RQ<sub>2</sub> *To what extent do influences of individual configuration options depend on the workload?*

#### B. Experiment Setup

1) *Subject System Selection*: We have selected nine configurable software systems for our study. To ensure that our findings are not specific to one domain or ecosystem, we include a mix of Java and C/C++ systems from different application domains (cf. Table I). In particular, we include systems studied in previous and related work [17], [35], [37], and we incorporate further ones with comparable size and configuration complexity (in terms of numbers of configurations

Table I: Subject System Characteristics

System	Lang.	Domain	Version	#O	#C	#W
JUMP3R	Java	Audio Encoder	1.0.4	16	4 196	6
KANZI	Java	File Compressor	1.9	24	4 112	9
DCONVERT	Java	Image Scaling	1.0.0- $\alpha$ 7	18	6 764	12
H2	Java	Database	1.4.200	16	1 954	8
BATIK	Java	SVG Rasterizer	1.14	10	1 919	11
XZ	C/C++	File Compressor	5.2.0	33	1 999	13
LRZIP	C/C++	File Compressor	0.651	11	190	13
x264	C/C++	Video Encoder	baee400...	25	3 113	9
z3	C/C++	SMT Solver	4.8.14	12	1 011	12

#O: No. of options, #C: No. of configurations, #W: No. of workloads tested

and configuration options). All systems operate by processing a domain-specific workload fed to them. Our study treats execution time as the key performance indicator with the exception of H2, for which we consider throughput.

2) *Workload Selection*: Our study relies on a selection of workloads for each domain or software system. Ideally, each set of workloads is diverse enough to be representative of most possible use cases. We selected the workload sets in this spirit, but cannot always guarantee a measurable degree of diversity and representativeness prior to conducting the actual measurements. Basically, this is what motivates this study in the first place. Nevertheless, we discuss this aspect as a threat to validity in Section VII.

Next, we outline the nine subject systems along with the workloads tested.

For the *audio encoder* JUMP3R, the measured task was to encode raw WAVE audio signals to MP3. We selected a number of different audio files from the Wikimedia Commons collection<sup>3</sup> and varied the file size/signal length, sampling rate, and number of channels. Both applications share all workloads.

For the *video encoder* x264, the measured task was to encode raw video frames (y4m format). We selected a number of files from the “derf collection”<sup>4</sup>, a set of test media for a variety of use cases. The frame files vary in resolution (low/SD up to 4K) and file size. For files with 4K resolution, we limited our measurements to encoding a subset of consecutive frames.

For the *file compression* tools KANZI, XZ, and LRZIP, we used a variety of community compression benchmarks that represent different goals, including mixes of files of different types (text, binary, structured data, etc.) or single-type files. We augmented this set of workloads with custom data, such as the Hubble Deepfield image and a binary of the Linux kernel. Beyond this set of workloads, for XZ and LRZIP, we added different parameterizations of the UIQ2 benchmark<sup>5</sup> to study the effect of varying file sizes.

For the *SMT solver* z3, the measured task was to decide the satisfiability (find a solution or counter example) of a range of logical problems expressed in the SMT2 format. We selected the six longest-running problem instances from z3’s performance test suite and augmented it with additional instances from the

<sup>3</sup><https://commons.wikimedia.org/wiki/Category:Images>

<sup>4</sup><https://media.xiph.org/video/derf/>

<sup>5</sup><http://mattmahoney.net/dc/uiq/>

SMT2-Lib repository<sup>6</sup>, to cover different types of logic and to increase diversity.

For the *SVG rasterizer* BATIK, the measured task was to transform a SVG vector graphic into a bitmap. We selected a number of resources from the Wikimedia Commons collection, primarily varying in terms of file size.

For the embedded *database* H2, we used a selection of four benchmarks (SmallBank, TPC-H, YCSB, Voter) from OLTPBENCH [47], a load generator for databases. For each benchmark, we varied the scale factor, which controls the complexity (number of entities modeled) in each scenario.

For the *image scaler* DCONVERT, the measured task was to transform resources (image files, Photoshop sketches) at different scales (useful for Android development). We selected files that reflect DCONVERT’s documented input formats (JPEG, PNG, PSD, and SVG) and vary in file size.

3) *Configuration Sampling*: For each subject system, we sampled a set of configurations. As exhaustive coverage of the configuration space is infeasible due to combinatorial explosion [48], for binary configuration options, we combine several coverage-based sampling strategies and uniform random sampling into an *ensemble* approach: We employ option-wise and negative option-wise sampling [2], where each option is enabled once (i.e., in, at least, one configuration), or all except one, respectively. In addition, we use pairwise sampling, where two-way combinations of configuration options are systematically selected. Interactions of higher degree could be found accordingly, however, full interaction coverage is computationally prohibitively expensive [48]. Last, we augment our sample set with a random sample that is, at least, the size of the coverage-based sample. To achieve a nearly uniform random sample, we used *distance-based sampling* [34]. If a software system exhibited numeric configuration options, we varied them across, at least, two levels to measure their effect.

4) *Coverage Profiling*: To assess what lines of code are executed for each combination of workload and software configuration, we used two separate approaches for Java and C/C++. For Java, we used the on-the-fly profiler JACOCO<sup>7</sup>, which intercepts byte code running on the JVM at run-time. For C/C++, we added instrumentation code to the software systems using CLANG/LLVM<sup>8</sup> to collect coverage information. We split the performance measurement and coverage analysis runs to avoid distortion from the profiling overhead.

5) *Measurements*: We conducted all experiments on three different compute clusters, where all machines within a compute cluster had an identical hardware setup: cluster A with an Intel Xeon E5-2630v4 CPU (2.2 GHz) and 256 GB of RAM, cluster B with an Intel Core i7-8559U CPU (2.7 GHz) and 32 GB of RAM, and cluster C with an Intel Core i5-8259U (2.3 GHz) and 32 GB of RAM. All clusters ran a headless Debian 10 installation (kernel 4.19.0-17 for cluster A and 4.19.0-14 for clusters B and C). To minimize measurement noise, we used

a controlled environment, where no additional user processes were running in the background, and no other than necessary packages were installed. We ran each subject system *exclusively* on a single cluster: H2 on cluster A; DCONVERT and BATIK on cluster B; the remaining systems on cluster C.

We collect performance data using the tools GNU TIME (execution time) and OLTPBENCH (throughput). For all data points, we report the median performance across five repetitions (except for H2), which has shown to be a good trade-off between variance and measurement effort [49]. Across these repetitions, most configurations exhibit only little variation (e.g., only a few seconds for whole-system benchmarks which run for several minutes): The ratio of configurations with a coefficient of variation (standard deviation divided by the mean) of less than 10 % ranges from 91 % (LRZIP) to 99 % (X264). For H2, we omitted the repetitions as, in a pre-study running on the identical cluster setup, we found that, across all benchmarks, the coefficient of variation of the throughput was consistently below 5 %.

## IV. STUDY RESULTS

In this section, we present the results of our empirical study with regard to variation of system-level performance distributions (RQ<sub>1</sub>) and the performance influence of individual configuration options (RQ<sub>2</sub>).

### A. Comparing Performance Distributions (RQ<sub>1</sub>)

1) *Operationalization*: We answer RQ<sub>1</sub> by pairwise comparing the performance distributions from different workloads (cf. the comparison in Figure 1) and by determining whether any two distributions are similar or, if not, can be transformed into each other: For the former case, we tested all combinations of workload-specific performance observations with the Wilcoxon signed-rank test<sup>9</sup> [52]. We rejected the null hypothesis  $H_0$  at  $\alpha = 0.95$ . To account for overpowering due to high and different sample sizes (cf. Table I), we further checked effect sizes to weed out negligible effects. Following the interpretation guidelines from Romano et al. [53], for no combination, Cliff’s  $\delta$  [54] exceeded a threshold effect size of  $|\delta| > 0.147$ . For the latter case, we are specifically interested in what *type* of transformation is necessary as this determines *how* complex a workload interacts with configuration options. Specifically, we categorize each pair of workloads with respect to the following aspects:

- 1) *Linear Correlation*: To test whether both performance distributions are shifted by a constant value or scaled by a constant factor, we compute for each pair of distributions Pearson’s correlation coefficient  $r$ . To discard the sign of relationship, we use the absolute value and a threshold of  $|r| > 0.6$  to indicate a linear relationship.
- 2) *Monotonous Correlation*: We test whether there is a monotonous relationship between the two performance distributions. We use Kendall’s rank correlation coefficient

<sup>6</sup><https://smt-comp.github.io/2017/benchmarks.html>

<sup>7</sup>JACOCO: <https://www.jacoco.org/jacoco/trunk/doc/>

<sup>8</sup>LLVM: <https://clang.llvm.org/docs/SourceBasedCodeCoverage.html>

<sup>9</sup>We use non-parametric methods since performance-distributions are often long-tailed and multi-modal [50], [51] and thus fail to meet requirements for parametric methods.

Table II: Three disjoint categories and criteria of relationships between pairs of workload-specific performance distributions.

Category		Criteria
<b>LT</b>	Linear transformation	$r^* \geq 0.6$
<b>XMT</b>	Monotonous transformation	$r^* < 0.6$ and $\tau^* \geq 0.6$
<b>NMT</b>	Non-monotonous transformation	(otherwise)

Table III: Frequency of each category (cf. Table II) for each software system studied and pairs of workloads.

System	$\Sigma_{\text{pairs}}$	<b>LT</b>		<b>XMT</b>		<b>NMT</b>	
		<i>abs</i>	<i>rel</i>	<i>abs</i>	<i>rel</i>	<i>abs</i>	<i>rel</i>
JUMP3R	15	15	100.0 %	0	0 %	0	0 %
KANZI	36	28	77.8 %	4	11.1 %	4	11.1 %
DCONVERT	66	29	43.9 %	0	0 %	37	56.1 %
H2	28	13	46.4 %	0	0 %	15	53.6 %
BATIK	55	28	50.9 %	8	14.6 %	19	34.6 %
XZ	78	65	83.3 %	1	1.3 %	12	15.4 %
LRZIP	78	57	73.0 %	13	16.7 %	8	10.3 %
X264	36	36	100 %	0	0 %	0	0 %
Z3	66	10	15.2 %	1	1.5 %	55	83.3 %

$\tau$  [55] and a threshold of  $|\tau| > 0.6$  for a monotonous relationship.

Based on these two correlation measures, we composed three categories that each pair of performance distributions can be categorized into. If both distributions exhibit a strong linear relationship, we classify them as linearly transformable (**LT**). If we observe a strong monotonous, but not a linear relationship, we classify such pairs as exclusively monotonously transformable into a separate category (**XMT**). Last, we have the pairs with a non-monotonous relationship (**NMT**). We summarize the category criteria as well as the category counts in Table III.

2) *Results*: We list the results of our classification in Table III. The observed range of relationships across the nine software systems exhibit no type that prevails across all software systems. All software systems, at least in part, exhibit performance distributions that can be transformed into one another using a linear transformation (**LT**). In particular, for JUMP3R and X264, we observe solely such behavior. The presence of linear transformations corroborates experimental insights from Jamshidi et al., who encoded differences between performance distributions using linear functions [21].

Exclusively monotonous transformations (**XMT**) are the exception and are found only in five out of the nine systems (KANZI, BATIK, XZ, LRZIP, Z3), twice with only one workload pair each (XZ and Z3). For all, except two systems (JUMP3R and X264), we observe non-monotonous relationships (**NMT**) with differing prevalence. For three systems (DCONVERT, H2, and Z3), the majority of transformations required is non-monotonous; for the other four systems (KANZI, BATIK, XZ, LRZIP), more than 10 % of workload pairs fall into this category.

**Summary (RQ<sub>1</sub>)**: Varying the workload causes a substantial amount of variation among performance distributions. Across workloads, we observed *mostly linear* (for six of the nine subject systems), but to a large extent, also *non-monotonous* relationships (for three of the nine subject systems).

## B. Workload Sensitivity of Individual Options (RQ<sub>2</sub>)

1) *Operationalization*: To address RQ<sub>2</sub>, we need to determine the configuration options' influence on performance and assess their variation across workloads.

*Explanatory Model*: To obtain accurate and interpretable performance influences per option, we learn an explanatory performance model based on the entire sample set using multiple linear regression [1], [2], [9]. Each variable in the linear model corresponds to an option, and each coefficient represents the corresponding option's influence on performance. We limit the set of independent variables to individual options (rather than including higher-order interactions) to be consistent with the feature location used for RQ<sub>3</sub>, where we determine option-specific, yet not interaction-specific code segments.

*Standardization*: To facilitate the comparison of regression coefficients across workloads, we follow common practice in machine learning and standardize our dependent variable by subtracting the population's mean performance and divide the result by the respective standard deviation. Henceforth, we refer to these standardized regression coefficients as *relative performance influences*. A beneficial side effect of standardization is that the observed variation of regression coefficients for each configuration option cannot be attributed to shifting or scaling effects (**LT**). This way, we can pin down the non-linear or explicitly non-monotonous effect that workloads may exercise on performance.

*Handling Multicollinearity*: Multicollinearity is a standard problem in statistics and arises when features are correlated [56]. It can, for instance, be caused from groups of mutually exclusive configuration options and result in distorted regression coefficients [1]. Although the model's prediction accuracy remains unaffected, we cannot trust and easily interpret the calculated coefficients. To mitigate this problem and, in particular, to ensure that the obtained performance influences remain interpretable, we follow best practices and remove specific configuration options from the sample that cause multicollinearity [1]. For the training step, we exclude all mandatory configuration options since these, by definition, cannot contribute to performance variation. In addition, for each group of mutually exclusive configuration options, we discard one group member when learning a model. These measures reduced multicollinearity to a negligible degree [57]. After these corrections, we observed no configuration options exceeding a variance inflation factor (indicating multicollinearity) of 5.

2) *Results*: Our results show a wide variety of degrees of workload sensitivity. Due to the size of our empirical study and space limitations, we selected three configuration options that showcase different characteristic traits of workload sensitivity that we observed. An exhaustive analysis for all configuration



options is illustrated in terms of an interactive dashboard provided as supplementary material. We strongly invite the interested reader to use this interactive dashboard to explore all distributions and results obtained in this study.

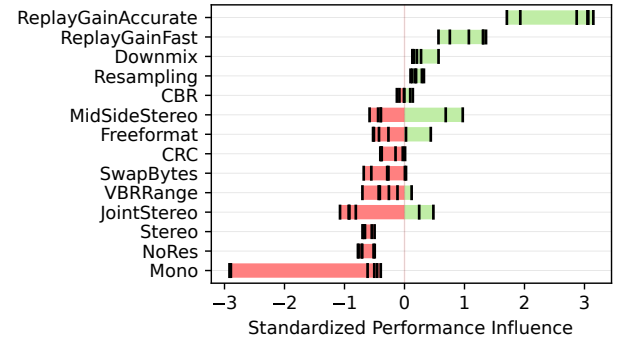
In Figure 2, we show the distribution of configuration options’ performance influences for three of the nine software systems (JUMP3R, Z3, and H2). Each vertical bar depicts the relative performance influence under a specific workload, the colored ranges depict positive (green) and negative (red) performance influences. The following patterns refer to one row in Figures 2a, 2b, and 2c (configuration option) for one subject system each.

*a) Conditional Influence:* For some configuration options, we observe that they affect performance only under specific workloads and remain non-influential otherwise. An example of such conditional influence is the configuration option **Mono** of the MP3 encoder JUMP3R. We illustrate the performance influence of this option across six workloads presented as bars in the last row of bars in Figure 3a. Selecting this option reduced the execution time substantially for two workloads, whereas, for the other workloads, the effect was far smaller.

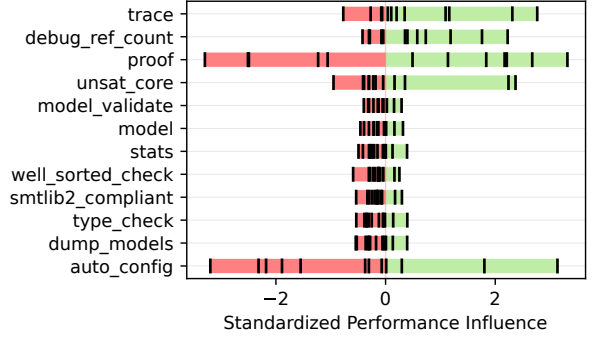
According to the documentation of JUMP3R<sup>10</sup>, selecting this option for stereo files (i.e., audio files with two channels) results in averaging them into one channel. Indeed, the two workloads described with reduced execution time the only ones that exhibit two audio channels in this selection. Hence, this example illustrates how a workload characteristic can condition the performance influence of a configuration option.

*b) High Spread:* Another pattern we found is that the performance influence of most (relevant) configuration options exhibits a large spread. For example, option **proof** of the SMT solver Z3 can both increase or decrease the execution time, as shown in Figure 3b (row 3). Compared to the example above, we cannot attribute this variation to an apparent workload characteristic. The global parameter **proof** enables tracking information, which is used for proof generation in the case a problem instance is unsatisfiable. Each workload in our selection contains multiple problem instances to decide satisfiability for. We conjecture that the ratio of satisfiable to unsatisfiable instances likely accounts for this high variation. From the user’s perspective, any input to the solver is opaque in that satisfiability as a workload characteristic cannot be determined practically without a solver.

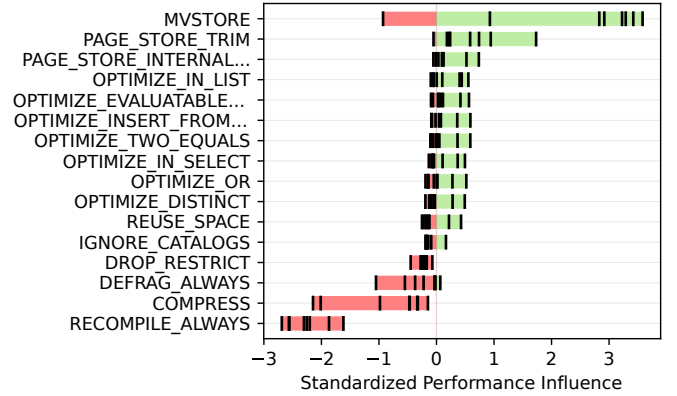
*c) Scaling Anomaly:* The scaling anomaly pattern is shown for the configuration option **MVSTORE** of the database system H2 in Figure 3c (left). This option controls which storage engine, either the newer multi-version store or the legacy page store, is used. We observe that selecting the newer multi-version store increases the measured throughput for all but one benchmark scenario. Using the Yahoo! Cloud Serving Benchmark (YCSB) with two different scale factors (which control the workload complexity, expressed as number of rows), we found that the lower complexity parameterization (ycsb-600) resulted in lower throughput. This is in stark contrast to



(a) Distribution of configuration options’ performance influences under different workloads for JUMP3R.



(b) Distribution of configuration options’ performance influences under different workloads for Z3.



(c) Distribution of configuration options’ performance influences under different workloads for H2.

Figure 2: Distribution of configuration options’ performance influences under different workloads for JUMP3R, Z3, and H2. Each vertical bar in a row depicts the performance influence of an individual option under a specific workload. The observed ranges of positive and negative influences are highlighted in green and red, respectively.

the expectation that a more complex workload would show lower throughput. While it is possible that some optimizations of the multi-version store are effective only under higher load, this example demonstrates that performance influence is not guaranteed to scale with the workload as expected.

<sup>10</sup>JUMP3R: <https://github.com/Sciss/jump3r/blob/master/README.md>

**Summary (RQ<sub>2</sub>):** Workloads can affect performance influences of configuration options in various ways (e.g., conditioning influence, introducing variance, having outliers). We can correlate some variation of performance influences with workload characteristics, yet identifying relevant workload characteristics is highly domain-specific and cannot be considered trivial.

## V. EXPLORING CAUSES OF WORKLOAD SENSITIVITY

The results for the first two research questions demonstrate sensitivity of the influence of configuration options to the workload due to non-monotonous interactions, which is consistent with the findings of a related study [45]. Before we discuss implications for performance modeling in the Section VI, we investigate the underlying factors that drive workload sensitivity.

We hypothesize that executions under different workloads also exhibit variation with respect to what code sections are executed (and which are not) and how this code is used (e.g., number of method invocations or loop iterations). Differences in performance influences of individual methods may stem from differences in program execution. Depending on whether an option is active or what value it has been set to, we may visit different code sections of the program or do so with varying frequency. To investigate to what extent one could infer or even explain performance variations based on different code execution profiles, we apply standard code-coverage analyses under different configurations and workloads. Our research question is as follows:

**RQ<sub>3</sub>** *Does the variation of performance influence of configuration options across workloads correlate with differences in the respective execution footprint?*

Exploiting a potential relationship between workload sensitivity of configuration options and differences in software execution could be beneficial. Instead of testing various combinations of configurations and workloads, code analysis can serve as a cost-effective way to detect workload sensitivity of individual options by identifying workload-specific differences in program execution.

### A. Operationalization

To explore whether and to what extent performance variations correlate with variations in the execution paths (that stem from the interplay of the given workload and configuration), we employ an analysis based on code coverage information (cf. Sec. III-B4).

We combine performance observations with code coverage data to evaluate the execution under different workloads, specifically focusing on code sections implementing option-specific functionality. By comparing the coverage of this code, we can develop hypothetical scenarios explaining workload sensitivity.

*First*, if we observe that the coverage of option-specific code is conditioned by the presence of some workload characteristic, we expect that such an option is influential only under the corresponding workloads. This scenario enables us (to some extent) to use code coverage as a cheap-to-compute proxy for estimating the representativeness of a workload and, by extension, resulting performance models: For options that are known to condition code sections, we can maximize option-code coverage to elicit all option-specific behavior and, thus, performance influence. For instance, a database system could cache a specific view only if a minimum number of queries are executed. Here, the effect of any caching option would be conditioned by the workload-specific number of transactions.

*Second*, if we observe performance variation across workloads in spite of similar or identical option-specific code coverage, we draw a different picture. In this case, we cannot attribute performance variation to code coverage, yet have to consider differences in the workloads' characteristics as potential cause: The presence of a workload characteristic may influence not *what* code sections are executed, but *how* code sections are executed. In a simple case, a software system's performance may scale linearly with the workload size. In a more complex case, the presence of a characteristic may determine how frequently an operation is repeated, as is the case for a table merge operation in a database system. Here, we would not elicit the worst-case performance if a previous transaction has sorted the data (e.g., by building an index).

1) *Locating Configuration-Dependent Code:* To reason about option-specific code, we require a mapping of configuration options to code. The problem of determining which code section implements which functionality in a software system is known as *feature location* [58]. While there is a number of approaches based on static [35], [59], [60] and dynamic analysis [36], [61], [62], we employ a more light-weight, but also less precise approach, that uses code coverage information. The rationale is that, by exercising feature code, for instance, by enabling configuration options or running corresponding tests, its location can be inferred from differences in code coverage. Applications of such an approach have been studied not only for feature location [63]–[66], but root in work on program comprehension [67]–[71] and fault localization [72], [73].

Specifically, we follow a strategy akin to *spectrum-based feature location* [65]: We commence with obtaining a baseline of all code that can be associated with a configuration option in the scope of our workload selection. Since we are looking for workload-specific differences in option-code coverage, the expressiveness of such a baseline depends on the diversity of the workloads in question. To infer option-specific code, we split our configuration sample (cf. Section III-B3) into two disjoint sets  $c_o$  and  $c_{-o}$  such that option  $o$  is selected only in  $c_o$  and not in  $c_{-o}$ . Next, we select from our code coverage logs the corresponding covered lines of code,  $S_o$  and  $S_{-o}$ . The rationale is that all shared lines between both sets are not affected by the selection of option  $o$ . Thus, we compute the *symmetric set difference*  $\mathbb{S}_o = S_o \Delta S_{-o}$  to approximate option-specific or,



at least, option-related code sections. To finally obtain code sections that are option-specific under a specific workload  $w$ , we repeat the steps above. Here, we consider only execution logs under workload  $w$  ( $S_{o,w}$  and  $S_{-o,w}$ ) and compute the symmetric set difference  $\mathbb{S}_{o,w} = S_{o,w} \Delta S_{-o,w}$ .

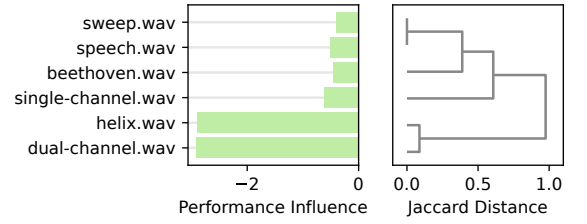
2) *Comparing Execution Traces:* From (a) the information about which code sections are specific to a configuration option and (b) how much of these sections is actually covered under different workloads, we can compare the workload-specific execution traces for each option. By comparing the sets  $\mathbb{S}_{o,v}$  and  $\mathbb{S}_{o,w}$  for any two workloads  $v$  and  $w$ , we can estimate the similarity between the option-code coverage with the Jaccard set similarity index. A Jaccard similarity of zero implies that there is no overlap in the code lines covered under each workload, whereas a Jaccard similarity of 1 implies that the exact same code was covered. Based on this pairwise similarity metric  $\text{sim}_o(v, w)$ , we can compute a corresponding distance metric  $d_o(v, w) = 1 - \text{sim}(v, w)$  and cluster all workload-specific execution profiles. We use agglomerative hierarchical clustering with full linkage to construct dendrograms. In this bottom-up approach, we iteratively add execution footprints to clusters and merge subclusters into larger ones depending on their Jaccard similarity to each other.

## B. Results

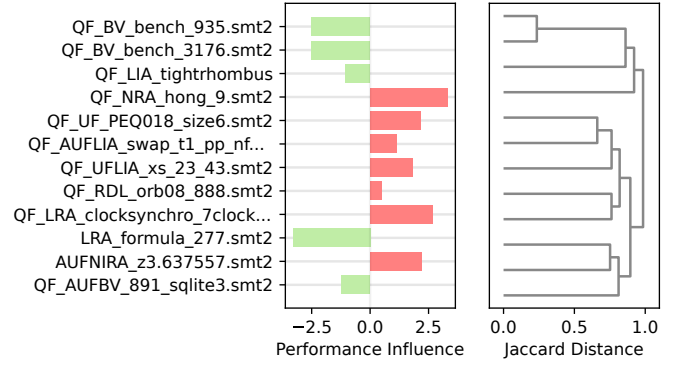
We report our findings for the relationship between execution footprints and performance influences for the same configuration options presented for  $RQ_2$ , since these illustrate likely causes of workload sensitivity and the limitations of solely relying on code coverage. The dashboard on the supplementary Web site provides diagrams and inferred feature code for all configuration options. The dendrograms next to the visualizations of performance influences in Figures 3a, 3c, and 3b, respectively, illustrate how similar the covered lines of option-specific code are under each workload. The dendrograms depict the Jaccard similarity clustering, where the split points indicate what Jaccard distance individual sets of lines or subclusters exhibit. The farther to the left the point is, the more similar are the components.

We observe that, in many cases, where a configuration option is “conditionally influential” (cf. Section IV-B2a), the respective option-specific code under the interacting workloads fall into a cluster, as with the option-specific code for Mono in Figure 3a. In this particular example, the dendrogram can be somewhat misleading as the number of common lines of code between workloads helix and dual-channel is far greater than between the other four workloads. Hence, differences in the coverage of option-specific code can account for, at least, some workload sensitivity.

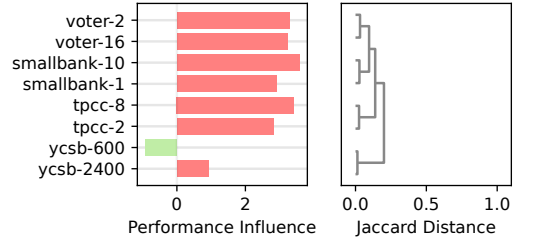
The other two examples, the configuration options **proof** (z3) and **MVSTORE** (H2), provide a different picture. Akin to the variation of performance influence of **proof**, the clustering for this configuration option (cf. Figure 3c) shows that some clusters are disjoint, and thus the option-specific code is highly fragmented depending on the workload.



(a) Workload-dependent performance influences of configuration option Mono of JUMP3R.



(b) Workload-dependent performance influences of configuration option proof of z3.



(c) Workload-dependent performance influences of configuration option MVSTORE of H2.

Figure 3: Workload-dependent performance influences of configuration options (a) Mono (JUMP3R), (b) proof (z3), and (c) MVSTORE (H2) (left) and option-code coverage clusterings for the the configuration options (right).

In the same vein, for **MVSTORE**, we see that the option-specific code is highly fragmented, yet all four benchmarks constitute clusters of high internal similarity. In the context of the observed variation of performance influences for the Yahoo! Cloud Serving Benchmark (YCSB), we see that even very high similarity in the covered code can virtually either improve or deteriorate performance.

For cases where we did not detect any differences in code coverage despite substantial differences in an option’s performance influence across workloads, our results suggest that the way *how* code was executed (i.e., how frequently methods or loops are executed) is shaping performance.

**Summary (RQ<sub>3</sub>):** Varying the workload can condition the execution of option-specific code coverage and cause performance differences. However, there is no single driving cause of variation: Code utilization depending on workload characteristics is a likely cause accounting for the majority of variation in the performance influence of an option.

## VI. DISCUSSION

Our results draw a clear picture of workload-induced performance variations of individual options. This sheds light on the extent of representativeness of single-workload performance models. But, this is not the end of the story: We saw complex variations that challenge transfer-learning approaches, which aim to overcome the workload specificity of models.

### A. Workload Sensitivity and Single-Workload Approaches

The observed workload sensitivity of configuration options exhibits a wide range of characteristics. While a large portion of options scales proportionally with workload complexity or remains unaffected by workload variation, the performance influences of several configuration options are sensitive to the workload. So far, the existing body of work on modeling [1]–[9] and optimizing [26]–[29] configuration-specific performance largely neglects the impact of workload variation at the cost of generalizability. Our findings from RQ<sub>2</sub> demonstrate that unexpected interactions of configuration options with the workload are not uncommon, which can distort performance estimations.

Beyond performance estimation, using performance models as surrogates for finding configurations with optimal performance properties is not without risk. For instance, there are several approaches utilizing the rank or importance of options [27], [29]. Given the observed workload sensitivity, such rankings remain susceptible to the choice of workload.

**Insight:** Workload sensitivity challenges the robustness and generalizability of single-workload performance models, yet it is neglected in state-of-the-art approaches. Worse, robust techniques using only rankings or relative importance of options are inapplicable for certain workload variations.

### B. Addressing Workload Variations

In Section II-C0a, we have outlined the existing body of work that aims at incorporating workload variations into performance modeling [19]–[22]. Despite the effectiveness of individual approaches, our results raise questions about assumptions used for transfer learning [20], [21] in this setting.

1) *Transfer Learning:* In their exploratory analysis, Jamshidi et al. reuse existing performance models by learning a linear transfer function across workloads [21]. Our results from RQ<sub>1</sub> have shown non-monotonous performance relationships across workloads, which is challenging to capture with such transfer functions. The presence of *non-monotonous interactions* between workloads and configuration options motivates employing more advanced machine learning techniques.

Table IV: Common top five influential configuration options among pairs of workloads.

System	Workload 1	Workload 2	# Common
JUMP3R	helix.wav	sweep.wav	2
KANZI	vmlinux	fannie_mae_500k	1
DCONVERT	jpeg-small	svg-large	2
H2	tpcc-2	tpcc-8	3
BATIK	village	cubus	4
XZ	deepfield	silesia	4
LRZIP	artifcl	uiq-32-bin	3
X264	sd_crew_cif_short	sd_city_4cif_short	4
Z3	QF_NRA_hong_9	QF_BV_bench_935	3

The more recent transfer-learning approach *Learning to Sample* [20] improves over the prior exploratory work by Jamshidi et al. [21]. It operates under the assumption that sampling for a new context, such as workloads, should focus on the influential options and interactions from a previously trained performance model. While this approach has shown to be effective, our results from RQ<sub>2</sub> contradict the basic assumption of stable influential options. To illustrate this in the context of our study, we select a pair of workloads for each of the nine subject systems studied and compare the ranking of configuration options with regard to their absolute performance influence (cf. RQ<sub>2</sub>). In Table IV, we show for each pair, how many configuration options are ranked in the top five (most influential) and shared across both workloads. For these workload pairs, we see that the rankings are inconsistent and thus not a reliable heuristic for transfer learning.

As the performance influence of configuration options can be conditioned by workload characteristics, a more appropriate metric to guide sampling would be to assess which configuration are workload sensitive rather than focusing on influential ones. This reiterates the problems described for most kinds of performance prediction approaches above.

2) *Workload-aware and Configuration-aware Performance Modeling:* While there is little work that *explicitly* considers the impact of factors beside configuration options on performance [19], our results from RQ<sub>2</sub> support idea of domain- or application-specific performance modeling. For instance, for several configuration options of JUMP3R, we can confidently associate workload sensitivity with a workload characteristic. To abstract more from application-specific approaches, a notion of workload sensitivity as a form of uncertainty is a promising avenue for further work. Work on using probabilistic programming to learn performance models [1] could be adapted to encode workload sensitivity.

**Insight:** Applying transfer learning to adapt performance models to new workloads must lift the assumption that the set of influential configuration options is stable. Domain-specific and workload-aware approaches are promising and should be extended on.

3) *Identifying Workload Sensitivity via Code Analysis:* Our findings from RQ<sub>3</sub> show that it is possible to identify workload

sensitivity through code analysis. This can be done using systematic coverage testing, which can be easily incorporated into CI/CD pipelines along with other code analyses, such as hit counts. While this low-cost metric can enhance existing approaches and help to interpret and contextualize performance estimations, it is important to note that more detailed analyses may be required to fully explain all performance variation. These findings can be applied in practice, for example, by using code coverage data to estimate up-front whether an option is input-sensitive and annotating existing performance models with a usage score per option. These results are important for understanding the performance of configurable software systems and for designing effective benchmarks.

**Insight:** Code analysis can be used to identify workload sensitivity and inform benchmark design in configurable software systems, but it is important to consider the limitations of this approach.

## VII. THREATS TO VALIDITY

Threats to *internal validity* include the presence of measurement noise, which may distort our classification into categories (Section IV-A) and model construction (Section IV-B). We address these threats by repeating each experiment five times (except for H2; cf. Section III-B5) and reporting the median as a robust measure in a controlled environment. Our coverage analysis (cf. Section III-B4) entails a noticeable instrumentation overhead, which may distort performance observations. We mitigate this threat by separating the experiment runs for coverage assessment and performance measurement. In the case of H2, the load generator of the OLTPBENCH framework [47] ran on the same machine as the database since we were testing an embedded scenario.

Threats to *external validity* include the selection of subject systems and workloads. To increase generalizability, we select software systems from various application domains as well as two different programming language ecosystems (cf. Table I). To increase the representativeness of our workloads, we vary relevant characteristics and, where possible, reuse workloads across subject systems of the same domain. Although there might be additional workload characteristics, our results demonstrate already for this selection severe consequences for existing performance modeling approaches. So, further variations could only strengthen our message.

## VIII. CONCLUSION

Configuration options are a key mechanism for optimizing the performance of modern software systems. Yet, state-of-the-art approaches of modeling configuration-dependent software performance largely ignore performance variation caused by changes in the workload. So far, there is no *systematic* assessment of whether, and if so, to what extent workload variations can render single-workload approaches inaccurate. We have conducted an empirical study of 25 258 configurations from nine configurable software systems to characterize the effects that varying workloads can have on configuration-specific

performance. We compare performance measurements with coverage data to identify possible similarities of executed code of different workloads compared to performance variations.

We find that workload variations affect software performance not only at the system-level, but also affect the influence of individual configuration options on performance, often in a non-monotonous way. While in some cases, we can correlate performance variations with the workload-conditioned execution of option-specific code, workload characteristics influence the utilization of option-specific code in further non-trivial ways (e.g., number of method calls).

We critically reflect on prevalent patterns, that we found in our subject systems and aim at raising awareness to the missing notion of workload sensitivity in existing approaches in this area. Our study provides an empirical basis that questions the practicality and generalizability of existing single-workload approaches as well as the validity of assumptions under which existing transfer-learning approaches in this area operate.

## IX. ACKNOWLEDGEMENTS

We thank our reviewers for their thoughtful and constructive comments. Apel’s work has been funded by the German Research Foundation (DFG) under contract AP 206/11-2 and grant 389792660 as part of TRR 248 – CPEC. Siegmund’s work has been supported by the German Research Foundation (DFG) under the contract SI 2171/2-2 and by the German Ministry of Education and Research (BMBF) and State Ministry for Science and Cultural Affairs of Saxony (SMWK) in the program Center of Excellence in AI research “Center for Scalable Data Analytics and Artificial Intelligence Dresden/Leipzig”, project identification ScaDS.AI.

## REFERENCES

- [1] J. Dorn, S. Apel, and N. Siegmund, “Mastering Uncertainty in Performance Estimations of Configurable Software Systems,” in *Proceedings of the International Conference on Automated Software Engineering (ASE)*. ACM, 2020, pp. 684–696.
- [2] N. Siegmund, A. Grebhorn, S. Apel, and C. Kästner, “Performance-Influence Models for Highly Configurable Systems,” in *Proceedings of the Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on the Foundations of Software Engineering (ESEC/FSE)*. ACM, 2015, pp. 284–294.
- [3] H. Ha and H. Zhang, “DeepPerf: Performance Prediction for Configurable Software with Deep Sparse Neural Network,” in *Proceedings of the International Conference on Software Engineering (ICSE)*. IEEE, 2019, pp. 1095–1106.
- [4] Y. Shu, Y. Sui, H. Zhang, and G. Xu, “Perf-AL: Performance Prediction for Configurable Software through Adversarial Learning,” in *Proceedings of the International Symposium on Empirical Software Engineering and Measurement (ESEM)*. ACM, 2020.
- [5] J. Guo, K. Czarnecki, S. Apel, N. Siegmund, and A. Wasowski, “Variability-aware Performance Prediction: A Statistical Learning Approach,” in *Proceedings of the International Conference on Automated Software Engineering (ASE)*. IEEE, 2013, pp. 301–311.
- [6] A. Sarkar, J. Guo, N. Siegmund, S. Apel, and K. Czarnecki, “Cost-Efficient Sampling for Performance Prediction of Configurable Systems,” in *Proceedings of the International Conference on Automated Software Engineering (ASE)*. IEEE, 2015, pp. 342–352.
- [7] J. Guo, D. Yang, N. Siegmund, S. Apel, A. Sarkar, P. Valov, K. Czarnecki, A. Wasowski, and H. Yu, “Data-efficient Performance Learning for Configurable Systems,” *Empirical Software Engineering*, vol. 23, no. 3, pp. 1826–1867, 2018.

- [8] Y. Zhang, J. Guo, E. Blais, and K. Czarnecki, "Performance Prediction of Configurable Software Systems by Fourier Learning," in *Proceedings of the International Conference on Automated Software Engineering (ASE)*. IEEE, 2015, pp. 365–373.
- [9] H. Ha and H. Zhang, "Performance-influence Model for Highly Configurable Software with Fourier Learning and Lasso Regression," in *Proceedings of the International Conference on Software Maintenance and Evolution (ICSME)*. IEEE, 2019, pp. 470–480.
- [10] J. Cheng, C. Gao, and Z. Zheng, "HINNPerf: Hierarchical Interaction Neural Network for Performance Prediction of Configurable Systems," *ACM Transactions on Software Engineering and Methodology (TOSEM)*, 2022.
- [11] S. Kounev, K.-D. Lange, and J. von Kistowski, *Systems Benchmarking*, 1st ed. Springer International Publishing, 2020.
- [12] S. Falkner, M. Lindauer, and F. Hutter, "SpySMAC: Automated Configuration and Performance Analysis of SAT Solvers," in *Theory and Applications of Satisfiability Testing (SAT)*, M. Heule and S. Weaver, Eds. Springer International Publishing, 2015, pp. 215–222.
- [13] L. Xu, F. Hutter, H. H. Hoos, and K. Leyton-Brown, "SATzilla: Portfolio-Based Algorithm Selection for SAT," *Journal of Artificial Intelligence Research*, vol. 32, no. 1, pp. 565–606, 2008.
- [14] Y. Ding, J. Ansel, K. Veeramachaneni, X. Shen, U.-M. O'Reilly, and S. Amarasinghe, "Autotuning Algorithmic Choice for Input Sensitivity," *SIGPLAN Not.*, vol. 50, no. 6, pp. 379–390, 2015.
- [15] D. Plotnikov, D. Melnik, M. Vardanyan, R. Buchatskiy, R. Zhuykov, and J.-H. Lee, "Automatic Tuning of Compiler Optimizations and Analysis of their Impact," *Procedia Computer Science*, vol. 18, pp. 1312–1321, 2013.
- [16] A. Maxiaguine, Y. Liu, S. Chakraborty, and W. T. Ooi, "Identifying 'Representative' Workloads in Designing MpSoC Platforms for Media Processing," in *Proceedings of the Workshop on Embedded Systems for Real-Time Multimedia (ESTImedia)*. IEEE, 2004, pp. 41–46.
- [17] J. A. Pereira, M. Acher, H. Martin, and J.-M. Jézéquel, "Sampling Effect on Performance Prediction of Configurable Systems: A Case Study," in *Proceedings of the International Conference on Performance Engineering (ICPE)*. ACM, 2020, pp. 277–288.
- [18] M. Khavari Tavana, Y. Sun, N. Bohm Agostini, and D. Kaeli, "Exploiting Adaptive Data Compression to Improve Performance and Energy-Efficiency of Compute Workloads in Multi-GPU Systems," in *International Parallel and Distributed Processing Symposium (IPDPS)*. IEEE, 2019, pp. 664–674.
- [19] U. Koc, A. Mordahl, S. Wei, J. S. Foster, and A. A. Porter, "SATune: A Study-Driven Auto-Tuning Approach for Configurable Software Verification Tools," in *Proceedings of the International Conference on Automated Software Engineering (ASE)*. IEEE, 2021, pp. 330–342.
- [20] P. Jamshidi, M. Velez, C. Kästner, and N. Siegmund, "Learning to Sample: Exploiting Similarities across Environments to Learn Performance Models for Configurable Systems," in *Proceedings of the Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on the Foundations of Software Engineering (ESEC/FSE)*. ACM, 2018, pp. 71–82.
- [21] P. Jamshidi, N. Siegmund, M. Velez, C. Kästner, A. Patel, and Y. Agarwal, "Transfer Learning for Performance Modeling of Configurable Systems: An Exploratory Analysis," in *Proceedings of the International Conference on Automated Software Engineering (ASE)*. IEEE, 2017, pp. 497–508.
- [22] P. Jamshidi, M. Velez, C. Kästner, N. Siegmund, and P. Kawthekar, "Transfer Learning for Improving Model Predictions in Highly Configurable Software," in *Proceedings of the International Symposium on Software Engineering for Adaptive and Self-Managing Systems (SEAMS)*. IEEE, 2017, pp. 31–41.
- [23] H. Martin, M. Acher, J. A. Pereira, L. Lesoil, J.-M. Jézéquel, and D. E. Khelladi, "Transfer Learning Across Variants and Versions: The Case of Linux Kernel Size," *Transactions on Software Engineering (TSE)*, vol. 48, no. 11, p. 42744290, 2022.
- [24] Y. Ding, A. Pervaiz, S. Krishnan, and H. Hoffmann, "Bayesian Learning for Hardware and Software Configuration Co-Optimization," University of Chicago, Tech. Rep. 13, 2020.
- [25] C. Lengauer, S. Apel, M. Bolten, A. Gröbinger, F. Hannig, H. Köstler, U. Rüde, J. Teich, A. Grebhahn, S. Kronawitter, S. Kuckuk, H. Rittich, and C. Schmitt, "ExaStencils: Advanced Stencil-Code Engineering," in *Parallel Processing Workshops – Euro-Par 2014 International Workshops, Revised Selected Papers, Part II*, ser. Lecture Notes in Computer Science, vol. 8806. Springer, 2014, pp. 553–564.
- [26] Chen, Tao and Li, Miqing, "Multi-Objectivizing Software Configuration Tuning," in *Proceedings of the Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on the Foundations of Software Engineering (ESEC/FSE)*. ACM, 2021, pp. 453–465.
- [27] V. Nair, T. Menzies, N. Siegmund, and S. Apel, "Using Bad Learners to Find Good Configurations," in *Proceedings of the Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on the Foundations of Software Engineering (ESEC/FSE)*. ACM, 2017, pp. 257–267.
- [28] V. Nair, Z. Yu, T. Menzies, N. Siegmund, and S. Apel, "Finding Faster Configurations Using FLASH," *Transactions on Software Engineering (TSE)*, vol. 46, no. 7, pp. 794–811, 2020.
- [29] J. Oh, D. Batory, M. Myers, and N. Siegmund, "Finding Near-Optimal Configurations in Product Lines by Random Sampling," in *Proceedings of the Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on the Foundations of Software Engineering (ESEC/FSE)*. ACM, 2017, pp. 61–71.
- [30] S. Apel, D. Batory, C. Kästner, and G. Saake, *Feature-Oriented Software Product Lines*. Springer, 2016.
- [31] C. Kaltenecker, A. Grebhahn, N. Siegmund, and S. Apel, "The Interplay of Sampling and Machine Learning for Software Performance Prediction," *IEEE Software*, vol. 37, no. 4, pp. 58–66, 2020.
- [32] F. Medeiros, C. Kästner, M. Ribeiro, R. Gheyi, and S. Apel, "A Comparison of 10 Sampling Algorithms for Configurable Systems," in *Proceedings of the International Conference on Software Engineering (ICSE)*. ACM, 2016, pp. 643–654.
- [33] N. Siegmund, S. Kolesnikov, C. Kästner, S. Apel, D. Batory, M. Rosenmüller, and G. Saake, "Predicting Performance via Automated Feature-Interaction Detection," in *Proceedings of the International Conference on Software Engineering (ICSE)*. IEEE, 2012, pp. 167–177.
- [34] C. Kaltenecker, A. Grebhahn, N. Siegmund, J. Guo, and S. Apel, "Distance-Based Sampling of Software Configuration Spaces," in *Proceedings of the International Conference on Software Engineering (ICSE)*. IEEE, 2019, pp. 1084–1094.
- [35] M. Velez, P. Jamshidi, F. Sattler, N. Siegmund, S. Apel, and C. Kästner, "ConfigCrusher: Towards White-Box Performance Analysis for Configurable Systems," *Automated Software Engineering (ASE)*, pp. 1–36, 2020.
- [36] M. Velez, P. Jamshidi, N. Siegmund, S. Apel, and C. Kästner, "White-Box Analysis over Machine Learning: Modeling Performance of Configurable Systems," in *Proceedings of the International Conference on Software Engineering (ICSE)*. IEEE, 2021, pp. 1072–1084.
- [37] M. Weber, S. Apel, and N. Siegmund, "White-Box Performance-Influence Models: A Profiling and Learning Approach," in *Proceedings of the International Conference on Software Engineering (ICSE)*. IEEE, 2021, pp. 1059–1071.
- [38] X. Han, T. Yu, and M. Pradel, "ConfProf: White-Box Performance Profiling of Configuration Options," in *Proceedings of the International Conference on Performance Engineering (ICPE)*. ACM, 2021, pp. 1–8.
- [39] S. Ceesay, Y. Lin, and A. Barker, "A Survey: Benchmarking and Performance Modelling of Data Intensive Applications," in *International Conference on Big Data Computing, Applications and Technologies (BDCAT)*. IEEE, 2020, pp. 67–76.
- [40] A. V. Papadopoulos, L. Versluis, A. Bauer, N. Herbst, J. v. Kistowski, A. Ali-Eldin, C. L. Abad, J. N. Amaral, P. Tüma, and A. Iosup, "Methodological Principles for Reproducible Performance Evaluation in Cloud Computing," *Transactions on Software Engineering (TSE)*, vol. 47, no. 8, pp. 1528–1543, 2021.
- [41] M. C. Calzarossa, L. Massari, and D. Tessera, "Workload Characterization: A Survey Revisited," *ACM Computer Survey*, vol. 48, no. 3, Feb. 2016.
- [42] Z. M. Jiang and A. E. Hassan, "A Survey on Load Testing of Large-Scale Software Systems," *Transactions on Software Engineering (TSE)*, vol. 41, no. 11, pp. 1091–1118, 2015.
- [43] L. Kotthoff, "Algorithm Selection for Combinatorial Search Problems: A Survey," in *Data Mining and Constraint Programming: Foundations of a Cross-Disciplinary Approach*, C. Bessiere, L. De Raedt, L. Kotthoff, S. Nijssen, B. O'Sullivan, and D. Pedreschi, Eds. Springer, 2016, pp. 149–190.
- [44] H. Samoaa and P. Leitner, "An Exploratory Study of the Impact of Parameterization on JMH Measurement Results in Open-Source Projects," in *Proceedings of the International Conference on Performance Engineering (ICPE)*. ACM, 2021, p. 213–224.

- [45] L. Lesoil, M. Acher, A. Blouin, and J.-M. Jézéquel, "The Interaction Between Inputs and Configurations Fed to Software Systems: An Empirical Study," *Computing Research Repository (CoRR)*, 2021.
- [46] P. Valov, J.-C. Petkovich, J. Guo, S. Fischmeister, and K. Czarnecki, "Transferring Performance Prediction Models Across Different Hardware Platforms," in *ICPE*. ACM, 2017, pp. 39–50.
- [47] D. E. Difallah, A. Pavlo, C. Curino, and P. Cudre-Mauroux, "OLTP-Bench: An Extensible Testbed for Benchmarking Relational Databases," in *Proceedings of the International Conference on Very Large Databases (VLDB)*, vol. 7, no. 4. VLDB Endowment, 2013, pp. 277–288.
- [48] C. Henard, M. Papadakis, M. Harman, and Y. Le Traon, "Combining Multi-Objective Search and Constraint Solving for Configuring Large Software Product Lines," in *Proceedings of the International Conference on Software Engineering (ICSE)*. IEEE, 2015, pp. 517–528.
- [49] I. Molyneaux, *The Art of Application Performance Testing*, 2nd ed., ser. Theory in Practice. Beijing: O'Reilly, 2015.
- [50] C. Curtsinger and E. D. Berger, "STABILIZER: Statistically Sound Performance Evaluation," in *Proceedings of the International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*. ACM, 2013, pp. 219–228.
- [51] A. Maricq, D. Duplyakin, I. Jimenez, C. Maltzahn, R. Stutsman, and R. Ricci, "Taming Performance Variability," in *USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, 2018, pp. 409–425.
- [52] M. Lovric, *International Encyclopedia of Statistical Science*, 1st ed. Springer, 2010.
- [53] J. Romano, J. D. Kromrey, J. Coraggio, J. Skowronek, and L. Devine, "Exploring Methods for Evaluating Group Differences on the NSSE and Other Surveys: Are the T-Test and Cohen's d Indices the Most Appropriate Choices?" in *Annual Meeting of the Southern Association for Institutional Research*, 2006, pp. 1–51.
- [54] N. Cliff, "Dominance statistics: Ordinal Analyses to Answer Ordinal Questions," *Psychological Bulletin*, vol. 114, pp. 494–509, 1993.
- [55] M. G. Kendall, "A New Measure of Rank Correlation," *Biometrika*, vol. 30, no. 1/2, pp. 81–93, 1938.
- [56] J. I. Daoud, "Multicollinearity and Regression Analysis," *Journal of Physics: Conference Series*, vol. 949, 2017.
- [57] R. M. O'Brien, "A Caution Regarding Rules of Thumb for Variance Inflation Factors," *Quality & Quantity*, vol. 41, no. 5, pp. 673–690, 2007.
- [58] J. Rubin and M. Chechik, "A Survey of Feature Location Techniques," in *Domain Engineering: Product Lines, Languages, and Conceptual Models*, I. Reinhartz-Berger, A. Sturm, T. Clark, S. Cohen, and J. Bettin, Eds. Springer, 2013, pp. 29–58.
- [59] M. Lillack, C. Kästner, and E. Bodden, "Tracking Load-Time Configuration Options," *Transactions on Software Engineering (TSE)*, vol. 44, no. 12, pp. 1269–1291, 2018.
- [60] L. Luo, E. Bodden, and J. Späth, "A Qualitative Analysis of Android Taint-Analysis Results," in *Proceedings of the International Conference on Automated Software Engineering (ASE)*. IEEE, 2019, pp. 102–114.
- [61] G. Bell, Jonatand Kaiser, "Phosphor: Illuminating Dynamic Data Flow in Commodity JVMs," *ACM SIGPLAN Notices*, vol. 49, no. 10, pp. 83–101, 2014.
- [62] C. H. P. Kim, D. Marinov, S. Khurshid, D. Batory, S. Souto, P. Barros, and M. D'Amorim, "SPLat: Lightweight Dynamic Analysis for Reducing Combinatorics in Testing Configurable Systems," in *Proceedings of the Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on the Foundations of Software Engineering (ESEC/FSE)*. ACM, 2013, pp. 257–267.
- [63] W. E. Wong and J. Li, "An Integrated Solution for Ttesting and Analyzing Java Applications in an Industrial Setting," in *Proceedings of the Asia-Pacific Conference on Software Engineering (ASPEC)*. IEEE, 2005, pp. 576–583.
- [64] M. Sulír and J. Porubán, "Semi-Automatic Concern Annotation using Differential Code Coverage," in *Proceedings of the IEEE International Scientific Conference on Informatics (ISCI)*. IEEE, 2015, pp. 258–262.
- [65] G. K. Michelon, B. Sotto-Mayor, J. Martinez, A. Arrieta, R. Abreu, and W. K. Assunção, "Spectrum-Based Feature Localization: A Case Study using ArgoUML," in *Proceedings of the International Conference on Software Product Lines (SPLC)*. ACM, 2021, pp. 126–130.
- [66] A. Perez and R. Abreu, "Framing Program Comprehension as Fault Localization," *Journal of Software: Evolution and Process*, vol. 28, pp. 840–862, 2016.
- [67] N. Wilde and C. Casey, "Early Field Experience with the Software Reconnaissance Technique for Program Comprehension," in *Proceedings of the International Conference on Software Maintenance (ICSM)*, 1996, pp. 312–318.
- [68] N. Wilde and M. C. Scully, "Software Reconnaissance: Mapping Program Features to Code," *Journal of Software Maintenance: Research and Practice*, vol. 7, no. 1, pp. 49–62, 1995.
- [69] K. D. Sherwood and G. C. Murphy, "Reducing Code Navigation Effort with Differential Code Coverage," Department of Computer Science, University of British Columbia, Tech. Rep. 14, 2008.
- [70] A. Perez and R. Abreu, "A Diagnosis-Based Approach to Software Comprehension," in *Proceedings of the International Conference on Program Comprehension (ICPC)*. ACM, 2014, pp. 37–47.
- [71] B. Castro, A. Perez, and R. Abreu, "Pangolin: An SFL-Based Toolset for Feature Localization," in *Proceedings of the International Conference on Automated Software Engineering (ASE)*. IEEE, 2019, pp. 1130–1133.
- [72] H. Agrawal, J. R. Horgan, S. London, and W. E. Wong, "Fault Localization using Execution Slices and Dataflow Tests," in *Proceedings of the International Symposium on Software Reliability Engineering (ISSRE)*. IEEE, 1995, pp. 143–151.
- [73] W. E. Wong, R. Gao, Y. Li, R. Abreu, and F. Wotawa, "A Survey on Software Fault Localization," *Transactions on Software Engineering (TSE)*, vol. 42, no. 8, pp. 707–740, 2016.