

# Is Code Coverage of Performance Tests Related to Source Code Features? An Empirical Study on Open-Source Systems

Muhammad Imran<sup>a</sup>, Vittorio Cortellessa<sup>a</sup>, Davide Di Ruscio<sup>a</sup>, Riccardo Rubei<sup>a</sup>, Luca Traini<sup>a</sup>

<sup>a</sup>*DISIM, University of L'Aquila, L'Aquila, 67100, Italy*

---

## Abstract

Performance testing aims to ensure the operational efficiency of software systems. However, many factors influencing the efficacy and adoption of performance tests in practice are not yet fully understood. For instance, while code coverage is widely regarded as a key quality metric for evaluating the efficacy of functional testing suites, there is limited knowledge about the types and levels of coverage that performance tests specifically achieve. Another important factor, often perceived as a barrier to the broader adoption of performance tests yet remaining relatively unexplored, is their extended execution time. In this paper, we examine (i) the coverage of performance testing suites, (ii) the characteristics of source code associated with performance-tested components, and (iii) the time cost of executing performance tests. Our analysis on 28 open-source systems reveals that performance tests achieve significantly lower code coverage than functional tests, as expected, and it highlights a significant trade-off between coverage and execution time. Our results also indicate a lack of generalizable characteristics in the source code covered by performance tests.

*Keywords:* Performance Testing, Unit Testing, JMH, JUnit, Java

---

## 1. Introduction

Software performance is a critical non-functional aspect of software systems. Deterioration in performance can present significant business challenges, including user dissatisfaction [1] and financial losses [2]. To address these concerns, organizations typically employ performance testing [3], a

technique designed to evaluate the software system’s performance before its deployment in a production environment. However, the practical implementation of performance testing faces challenges.

One of such challenges is associated with the development and maintainability of performance testing suites. Creating these tests often requires specific technical expertise that may not be readily available to the average developer [4, 5]. Additionally, software development processes tend to prioritize functional development activities, which can lead to limited resource allocation for quality assurance tasks [6, 4], such as the creation and maintenance of performance tests [5]. These factors collectively contribute to the oversight of performance assurance activities, leading to potential implications on the *quality* of performance testing suites.

The traditional way of assessing the *quality* of a testing suite involves using *code coverage*. This metric gauges the extent to which the testing suite executes the software source code, serving as a simple yet effective indicator of the testing suite quality. Although the validity of code coverage as a measure of test quality is still a matter of debate [7, 8], it remains the de-facto standard for evaluating testing suites in practical scenarios. Code coverage has traditionally been used and studied in the context of functional testing [9, 10, 7, 8] (e.g., unit tests), and there are increasing indications that its relevance extends to software performance testing as well. For instance, researchers have shown that code coverage significantly impacts the test capability of triggering performance bugs [11], and that extending coverage in performance testing suites can enhance their effectiveness in discovering performance issues [12].

Despite these indications, there is still limited knowledge about the code coverage achieved by performance tests. In a previous conference paper [13], we took an initial step on this issue by conducting an investigation into the code coverage of performance tests and the associated execution time cost, which is an important factor typically in trade-offs with code coverage [14]. Our methodology involved selecting 28 Java software systems on GitHub that featured JMH benchmarks, a widely used form of performance tests in Java. We then conducted a comprehensive analysis to extract all Java methods within these software systems. Finally, we executed 2,190 JMH benchmarks on these Java software systems to identify the methods covered by performance testing. Our findings revealed that JMH benchmarks achieve a limited code coverage of about 8.8% on average, which is 4 times lower than the JUnit tests one. Additionally, JMH benchmarks incur a considerably high

execution time cost, which is in average 62 times larger than those of JUnit tests.

This paper extends our previous one [13], VITTORIO ▶with the main intent to investigate whether static features can play the role of “detectors” of code that is usually covered by performance tests. The identification of one or more such features, in perspective, would support the work of performance testers by pointing them to portions of code that claim their attention, only on the basis of static code analysis, therefore without needing to execute those portions of code. With this main motivation in mind,◀ we introduced a wider analysis that explores the correlation between the characteristics of code components and their likelihood of being covered by performance tests. Specifically, we aim to understand whether code components covered by performance tests share characteristics (e.g., recurrent patterns) that are generalizable across software systems. The identification of such patterns can help, in particular: (i) to guide developers in selecting which code components to test, and (ii) to improve performance engineering techniques, such as performance test prioritization [15, 16], selection [17], and generation [12, 18].

In order to investigate this aspect, we extracted 44 well-known source code metrics from 177, 697 Java methods, and we trained machine learning classifiers to predict whether a method is covered by a performance test based on these metrics. Through this process, we searched for correlations between the code metrics and the likelihood of being covered by performance tests. We employed a rigorous methodology that involved training 11 machine learning algorithms across 99 configurations, while adhering to best practices such as dataset rebalancing, feature selection, and hyper-parameter tuning. Our results basically exclude correlations between source code metrics and performance test coverage, by suggesting a lack of generalizable recurrent patterns in performance-tested methods. This finding indicates that the choice of code components to be performance tested is primarily driven by domain-specific objectives rather than generalizable code characteristics. Consequently, this poses significant challenges for the automated generation, selection, and prioritization of performance tests.

To sum up, this paper provides the following contributions:

- An empirical investigation on the code coverage of performance testing, and the associated execution time cost.
- An extensive investigation on the correlation between source code metrics and performance testing coverage.

- A replication package [19] containing the dataset we mined in the study, the scripts we used to perform the data analysis, and the detailed results of our analysis.

## 2. Background

### 2.1. Performance Testing

Performance indicates how well a software system or component operates in terms of efficiency, including aspects like response time and throughput. The *response time* refers to the time required to respond to a request, while the *throughput* of a system is the number of requests that can be processed in some specified time interval [20]. Just like functional correctness, performance is one of the key characteristics in different classifications of the software quality attributes [21][22].

A common practice to determine the performance of any application is performance testing [23]. Performance testing is an activity aimed at evaluating the efficiency of a software system in a pre-production environment, often in terms of response time, throughput, or combination of both [24]. It includes a wide variety of approaches, such as system-level testing and load testing [25]. One of the goal of performance testing is to identify problems that cause the system performance to degrade. Even if a software system lacks explicit performance specifications, it is implicitly required that it should not take infinite time or resources to execute.

Performance testing is usually carried out by repetitive executions of target code until enough performance data points are collected [26], [27]. To evaluate overall software systems, state-of-the-art performance testing practices typically involve relatively large-scale and long-running tests [28]. However, this process should be systematic due the the expensive nature of performance tests in terms of execution time cost [13]. To deal with this intrinsic complexity, several approaches to performance unit testing, such as microbenchmarking, have been proposed for accurate performance evaluation of small-scale and isolated source code segments [12].

Microbenchmarks are small test cases used to evaluate the performance of a specific piece of code or function [29]. They are increasingly becoming popular and widely adopted, mainly due to their short runtimes and testing durations [30]. Performance microbenchmarking aims to test smaller code units (e.g. methods or statements) [31], [32]. This approach is particularly

important in languages like Java, where system performance can be significantly affected by apparently small code segments. Microbenchmarking is used to evaluate specific performance metrics such as execution time, memory consumption, and throughput, while providing detailed insights on the performance characteristics of individual code units.

However, there are inherent challenges to performance microbenchmarking, like unreliable results [33], and a lack of appropriate tools [34]. Due to these challenges, many frameworks provide benchmark tooling like Caliper [35], AutoJMH [18], and JUnitPerf [36]. For Java-based software, JMH [37] is the de facto standard framework used to automate performance microbenchmarking by defining and executing software benchmarks. Therefore, in this study, we focus on microbenchmarking using JMH in Java programs where we consider methods as units to be tested.

## *2.2. Test Coverage*

Test coverage is a commonly used metric in software testing that calculates how thoroughly the test cases exercise a given program. It is used as an indicator to monitor the quality of a test suite and a way to measure how thoroughly the software is tested. Test coverage refers to the portion of a software application used during a particular test execution [38]. Test coverage also serves to inform analysis techniques, such as fault localization, or to guide search-based algorithms towards generating coverage-optimized test suites (e.g., [39], [40]).

Typically, test coverage criteria include statement coverage and branch coverage at the level of methods, classes, packages, and overall system. Statement coverage measures the percentage of executable statements exercised by a test suite, while branch coverage measures the percentage of branches exercised by a test suite. The coverage can be determined using various test coverage tools. The coverage measurement is performed by inserting measurement probes into a program at the targeted granularity level, such as at a statement, block, or method level, and then running the tests on the instrumented program, while typically storing coverage information in trace files [41]. All the tools usually work on this principle but they may vary in the languages to which they apply. Also, they may provide measurements at different granularity levels such as statement, branch, method, and class [42].

Despite limited applications and attention to test coverage for performance testing, better test coverage has proven to be effective in finding syn-

chronization performance bugs in production environments [43]. Applying test case prioritization, coverage-based techniques provide significant benefits in regression testing by efficiently uncovering performance regressions [44]. Studies have shown that dynamic coverage analysis (e.g., line coverage or branch coverage that involve measuring the test coverage while the program is running) can effectively measure the execution behavior of Java programs. These techniques prove to be effective in early capturing significant performance changes [45]. Due to this, coverage-based approaches are particularly beneficial for performance testing of software microbenchmarks, which typically require longer execution times than unit tests [13].

### 3. Study Design

This study investigates the code coverage achieved by *performance tests* and their time cost, given the relevance of this aspect for the practical usage of performance tests. To aid the interpretation of our results, we conduct a comparative analysis with *functional tests* to assess the differences both in terms of code coverage and time cost. Furthermore, we explore the potential relationships between source code metrics and the likelihood of the analyzed source code being covered by performance tests. By identifying these relationships, we aim to uncover recurring patterns in the source code that can inform and guide performance testing efforts.

We focus on JMH microbenchmarks and JUnit tests due to their extensive use within the Java ecosystem. JMH is the de-facto standard for developing and running microbenchmarks, a widely known form of performance tests in Java software [46], while JUnit stands as one of the most popular libraries for implementing and executing Java functional tests. In this study, we aim to address the following research questions (RQs):

▷ **RQ<sub>1</sub>**: *To what extent do performance tests cover the source code of software systems?* Our objective is to evaluate the extent of code coverage achieved by performance testing suites across various software systems and compare it to the coverage provided by functional testing suites. We assess coverage from multiple perspectives, including overall and direct coverage, and we analyze the degree of overlap among different tests in their coverage of identical sections of the source code.

▷ **RQ<sub>2</sub>**: *What is the time cost of performance tests?* We assess the time costs incurred during the execution of performance tests, and compare them

to those associated with functional tests. We measure the overall time consumed at *suite-level* and *test-level*, thus providing insights on the temporal impact of performance testing at different granularity.

▷ **RQ<sub>3</sub>**: *Are there any specific characteristics in the source code that can guide performance testing?* To determine if there are any peculiarities in the source code sections chosen for performance testing, we investigate the relationship between static code features and the likelihood of being covered by performance tests. To do so, we leverage machine learning models to predict whether a method is performance tested or not. Specifically, we employ a carefully designed process consisting of different steps, while adhering to current machine learning best practices largely inspired from a similar work [30]. We start with machine learning models using default settings. Then, we apply both class rebalancing techniques [47, 48] and feature selections [49, 50]. Finally, we adopt hyperparameter tuning [51] to enhance the prediction accuracy of the models. High prediction accuracy indicates the presence of distinctive patterns in the source code that are subject to performance testing. Conversely, low prediction accuracy suggests a lack of generalized characteristics in performance tested methods. To address **RQ<sub>3</sub>**, we investigate the following three sub-research questions:

- **RQ<sub>3.A</sub>**: *Can we predict whether a method should be performance tested using machine learning models with default settings?* We investigate the use of machine learning models with their default settings to evaluate if static code features can be used to predict whether a code component should be covered by performance tests.
- **RQ<sub>3.B</sub>**: *How does the prediction accuracy change when employing feature selection and class-rebalancing techniques?* We employ two well-known preprocessing steps, namely (i) feature selection and (ii) class-rebalancing, to investigate their impact on the prediction accuracy of machine learning models. For feature selection, we remove less relevant code features, such as co-linear and multi-co-linear features, from our dataset. For class-rebalancing, we apply both oversampling and undersampling techniques.
- **RQ<sub>3.C</sub>**: *How does the prediction accuracy change after tuning the hyperparameters of machine learning models?* We aim to investigate the extent to which the optimization of hyperparameter configurations can enhance the prediction accuracy of machine learning models.

### 3.1. Main steps of the performed study

Given the objective of our study, we selected an initial pool of 40 Java software projects hosted on GitHub. This selection was guided by four key considerations: (i) these systems are well-established Java libraries that cover a broad spectrum of application domains, (ii) each of these projects includes JMH benchmarks and JUnit tests, (iii) we are familiar with the commands necessary to execute the JMH microbenchmarks for these systems, and (iv) they have been used in prior work [30, 16, 52, 53, 54], thus supporting their appropriateness in this context.

As shown in Figure 1, the executed process consists of three main steps: (i) *raw data collection* from the selected software systems, (ii) *data wrangling*, and (iii) *model training and tuning* phase, as described in the following subsections.

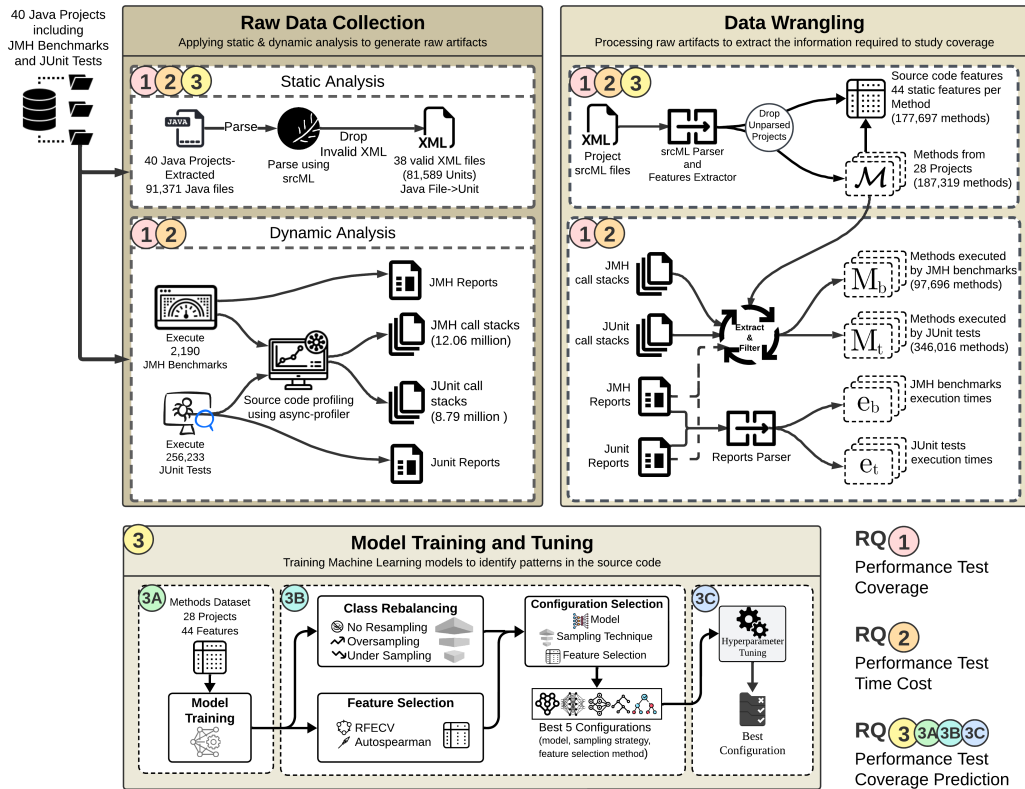


Figure 1: Main steps of the performed study.

### 3.1.1. Raw Data Collection

Our data collection combines both static code analysis and dynamic analysis of test executions. In particular, we collected three distinct types of raw data from each Java system, i.e., *srcML XML files*, *JMH/JUnit callstacks*, and *JMH/JUnit reports*.

▷ *srcML XML files*: We transformed each Java source file into a structured XML format using srcML toolkit [55]. This transformation facilitates a structured and reliable static analysis of the source code, thus enabling straightforward extraction of relevant code information, such as the list of Java method signatures and other source code elements corresponding to Java language features for specifying and managing control flows, data, and concurrency that appear in the project. In this process, we encountered a specific issue with two projects, namely *apache hive* and *eclipse jersey*, where the srcML toolkit generated incomplete XML files that omitted the representation of specific Java source files. Thus, we decided to exclude these two projects from our analysis. As a result, we successfully parsed the source code of 38 of our initially selected projects.

▷ *JMH/JUnit callstacks*: We rely on dynamic analysis to identify the code components covered (or not covered) by JMH/JUnit tests. We prefer dynamic analysis over static analysis due to several limitations of the latter [56], such as its inability to reliably derive method invocations. Specifically, we employed *async-profiler*<sup>1</sup> to profile the execution of JMH benchmarks and JUnit tests by capturing their respective call stacks. A call stack reports the currently active methods in the CPU and the sequence of their invocations. For illustration, consider Fig. 2, which presents a call stack from the execution of a JMH benchmark named `skinnyEncodeIntoCompressedByteBuffer`. Through this call stack, we can deduce the sequence of executed methods within the benchmark. From Fig. 2, we can observe that the benchmark first calls `encodeIntoCompressedByteBuffer`, which subsequently invokes `encodeIntoByteBuffer`, and this is followed by `writeCountsDiffs`, and so forth.

After executing each testing suite, we produced a distinct file that catalogs all unique call stacks observed during the execution. However, during this procedure, we faced issues attaching the profiler to the JMH benchmarks of seven projects. Consequently, we excluded these projects from our anal-

---

<sup>1</sup><https://github.com/asynct-profiler/async-profiler>

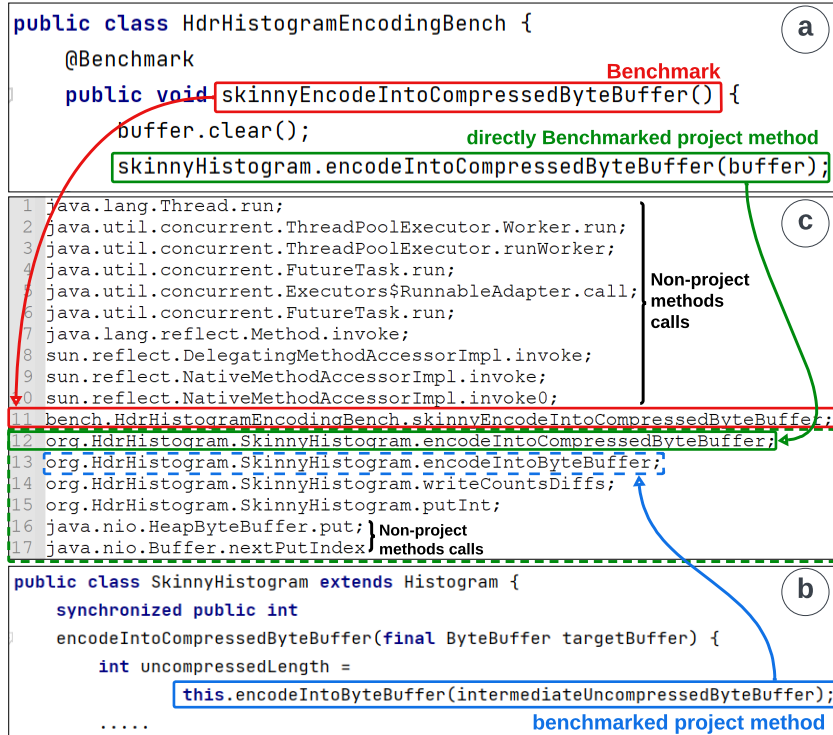


Figure 2: (a) Example invocation of a method from a benchmark. (b) Benchmarked method definition and indirect call to another project method. (c) Example of a collected call stack.

ysis. We faced similar issues for JUnit testing in nine projects. As a result, we collected 12.06 million unique call stacks for JMH benchmarks from 31 systems, and 8.79 million unique call stacks for JUnit tests from 19 projects.

▷ *JMH/JUnit reports*: We collected JMH and JUnit reports to obtain two primary information: (i) the list of JMH and JUnit tests and (ii) their corresponding execution times. For the latter metrics, we have exploited JMH configuration mechanisms, which allow developers to declare the number of repetitions for each JMH test, thus to deduce a bound on the time needed to execute the JMH test. To obtain the JMH configurations for each JMH test, we applied an approach similar to a prior work [54], which exploits a JMH feature that allows to overwrite configurations on the fly via CLI arguments. We executed each JMH test while reducing the execution time through JMH

CLI arguments<sup>2</sup>, and we stored the associated JMH reports, which include the JMH configurations set by developers. To gather execution times of JUnit tests, we utilized the Maven Surefire plugin<sup>3</sup>. This plugin is an established instrument within the Java ecosystem, and it is explicitly tailored for running JUnit tests through Maven. The execution of tests produces XML reports that break down the execution time for each JUnit test.

We conducted all tests on a dedicated machine equipped with Linux Ubuntu 18.04.2 LTS, powered by a dual Intel Xeon CPU E5-2650 v3 at 2.30 GHz, boasting 40 cores and 80 GB of RAM.

### 3.1.2. Data Wrangling

To address our research questions, we focus on four key pieces of information for each Java system: (i) the entire set  $\mathcal{M}$  of Java methods appearing in source code, (ii) the set of features for each method (as listed in Table 1), (iii) the set  $M_t$  of methods covered by each (JMH/JUnit) test  $t$ , and (iv) the execution time  $e_t$  of each (JMH/JUnit) test  $t$ . In the following, we describe the process used to derive this information, starting from the raw data.

▷ *Java methods*: To extract the fully qualified names of all methods within each project, we parsed the *srcml XML files* using the *lxml* library and employing XPath queries. While executing XPath queries on XML files, facilitated by *lxml*<sup>4</sup>, we encountered a limitation regarding the size capacity. Specifically, there was a size threshold (i.e., 6,800 srcML units) over which *lxml* could not evaluate the given XPath expressions. This constraint necessitated the exclusion of three projects from our analysis. Ultimately, this process resulted in extracting a set  $\mathcal{M}$  of Java methods for each Java system. In total, we extracted 187,319 methods from 28 distinct systems (as reported in Table 5, where detailed information about the amount of *Methods*, *Benchmarks* and *Unit Tests* per project is provided).

▷ *Source code features*: To extract the features from the source code, we employed additional XPath queries on the *srcml XML files*. Specifically, for each Java method in  $\mathcal{M}$ , we extracted a set  $\mathcal{F}$  of 44 source code features. The complete list of considered features is available in Table 1. To maintain consistency, we have adhered to specific naming conventions for the employed features. In particular, all features are named using Camel case notation.

---

<sup>2</sup>We refer the reader to [54] for a detailed explanation of this process.

<sup>3</sup><https://maven.apache.org/surefire/maven-surefire-plugin/>

<sup>4</sup><https://lxml.de/>

Table 1: Collected source code features.

Category/Level	Feature Name	Description	
(i) meta information	class	classScope	access specifier of the class
		cIsAbstract	if the class is declared abstract
		cNestingLevel	nesting depth of the class (0 for no nesting)
		#cImports	no. of import statements
		#cNativeImports	no. of native imports
		#cMethods	no. of methods
		#cInherits	if the class uses inheritance
		#cImplements	no. of interfaces implemented
		#classLOC	Lines of code of the class
		#cPkgClasses	no. of classes in current Package
	method	#methodLOC	Lines of code of the method
		#nameLen	method name length
		methodScope	access specifier of the method
		isOverloaded	if the method is overloaded
(ii) language feature	control flow and data	#if	no. of if conditions
		#switch	no. of switch statements
		#case	no. of case statements
		#for	no. of for loops
		#while	no. of while loops
		#do	no. of doWhile loops
		#nestedLoops	no. of nested loops (arbitrary depth)
		#methodCalls	no. of methods called
		#internalCalls	no. of internal methods called
		#externalCalls	no. of external methods called
		#return	no. of return statements
		#throw	no. of throw statements
		#catch	no. of catch statements
#cyclo	cyclomatic complexity by McCabe [57]		
#vars	no. of the variables declared		
(iii) calls to standard library APIs	calls to methods from selected standard libraries in Java (details are given in Table 2)		

Moreover, quantitative feature names begin with the prefix “#”, whereas features at the *class granularity level* include the prefix “c”.

We considered three kinds of source code features: (i) *meta information*, e.g., the number of LOCs or files (these features were collected at both class and method granularity level); (ii) *language features*, e.g., loops, conditionals, or variables; and (iii) *calls to standard library APIs*. We have chosen these kinds of source features because of their relationship with software performance, which might influence the likelihood of the Java method being performance tested or not [58]. The last reference supports the last sentence of the para (left as is). Previous research on software performance often leverages statically or dynamically determined source code features, such as added loops or method calls, to predict whether a code change slows down (or speeds up) the program under analysis [59].

Below, we provide a more detailed description of each kind of source code features.

- (i) **Meta Information:** Meta information features provide insights into the structure and characteristics of the code, which can potentially impact the software performance. We collected these features at both the class and method granularity levels. At the class level, features such as the number of import statements (#cImports), number of methods (#cMethods), and number of lines of code in a class (#class-LOC) offer a high-level view of the codebase. At the method level, features such as the number of lines of code in a method (#method-LOC), the method name length (#nameLen), and the access specifier of the method (methodScope) provide more fine-grained information. Research has shown that meta information features, including file and method granularity features, have been successfully employed in prediction models for performance properties [58].
- (ii) **Language Features:** They are related to control flow and data elements, which have been extensively used in prior work [60, 61, 58]. Control flow elements include language constructs to specify and manage conditions (#if, #switch, #case), loops (#for, #while, #do, #nestedLoops), function lifecycle (#methodCalls, #internalCalls, #externalCalls, #return), exception handling (#throw, #catch), and cyclo-matic complexity (#cyclo). These elements contribute to the increased complexity of a function, which may negatively impact software performance. Research has frequently identified loops as a root cause

of performance issues [62, 63, 64, 65], and control flow elements have been employed in machine-learning-based performance prediction models [58, 11]. Also, the method memory usage and stack size may potentially affect performance due to increased garbage collection (GC) activity. So, more variables (`#vars`) increase pressure on the GC, thus leading to more frequent memory allocations and deallocations, which can impact performance [66].

- (iii) **Library Calls:** Standard library packages provide functionalities such as file and network I/O, communication with the operating system (OS), text and string processing, and concurrency primitives. These functionalities often involve non-deterministic behaviors, such as waiting for locks, blocking during I/O operations, and network data transmission, which can affect software performance. Consequently, calls to certain standard libraries may adversely impact software performance. Additionally, inefficient or incorrect use of APIs, as well as concurrency and synchronization issues, have been identified as common sources of performance bugs [62, 63, 64, 65]. To analyze standard library calls, we group them as single features at the package level. All calls to a specific package are combined into one feature. Table 2 shows the Java standard packages mapped as boolean features with their descriptions.

▷ *Test coverage:* We leveraged the JMH and JUnit call stacks to identify the Java methods covered by each JMH benchmark or JUnit test. For each test  $t$ , we extracted the set  $M_t$  of project methods executed within  $t$ . To achieve this, we iterated over all the project methods in  $\mathcal{M}$  and checked if each method appeared after  $t$  in at least one call stack. A method  $m$  is considered covered by test  $t$  if it is invoked either directly or indirectly by  $t$  (i.e.,  $m \in M_t$ ). Additionally, for each test  $t$ , we created a separate set  $\hat{M}_t$  that only contains the methods directly called by  $t$ . In other words, for each test  $t$ , we select the methods in  $\mathcal{M}$  that appear immediately after  $t$  in the call stacks. For instance, in the example shown in Fig. 2, the method directly invoked by the benchmark `skinnyEncodeIntoCompressedByteBuffer` would be `encodeIntoCompressedByteBuffer`, but not `encodeIntoByteBuffer`. At the end of this process, for each JMH benchmark or JUnit test  $t$ , we obtain two sets:  $M_t$  representing methods covered either directly or indirectly by  $t$ , and  $\hat{M}_t$  representing methods directly covered by  $t$ .

▷ *Test execution time:* The JMH configuration set by developers determines the execution time of a JMH benchmark. This configuration defines

Table 2: Standard library call features from Table 1

Package	Feature Name	Category	Description
java.util	usesJavaUtil	utility	Utility classes in Java.
java.lang	usesJavaLangThread	concurrency	Classes for multi-threading and concurrent programming.
java.util.concurrent	usesJavaUtilConcurrent	concurrency	Advanced concurrency utilities & collections.
java.io	usesJavaIo	io	Classes for I/O through data streams.
java.nio	usesJavaNio	io	New I/O for more scalable I/O operations.
java.nio.channels	usesJavaNioChannels	io	Channels and selectors for non-blocking I/O.
java.nio.file	usesJavaNioFile	io	File I/O enhancements using the NIO framework.
java.nio.charset	usesJavaNioCharset	io	Charset classes for encoding and decoding.
java.net	usesJavaNet	io	Networking classes for implementing protocols and communication.
javax.net.ssl	usesJavaxNetSsl	io	SSL/TLS support for secure network communication.
java.lang	usesJavaLang	os	Core Java language classes and basic types.
java.lang.management	usesJavaLangManagement	os	Management interfaces for the Java platform.
java.util.regex	usesJavaUtilRegex	strings	Regular expressions for pattern matching and string manipulation.
java.text	usesJavaText	strings	Classes for parsing and formatting text.
java.math	usesJavaMath	math	Mathematical functions and utilities.

the levels of repetitions (i.e., forks, iterations, and invocations) used during benchmarking to address the inherent variability of performance measurements [67]. Invocations are repeated benchmark executions within a time-bound iteration, while a series of iterations forms a fork. Each fork usually comprises two distinct types of iterations: warmup and measurement iterations. Warmup iterations are intended to bring the fork into a steady state of performance [54, 67], while measurement iterations are the ones that are actually used for performance assessment. For each benchmark  $t$ , we first extract the JMH configuration from the JMH reports, i.e., the warmup iteration time  $w$ , the measurement iteration time  $r$ , the number of warmup iterations  $wi$ , the number of measurement iterations  $i$ , and the number of forks  $f$ . Then, we compute the associated execution time  $e_t$  accordingly:

$$e_t = (w \cdot wi + r \cdot i) \cdot f.$$

For JUnit tests, we instead directly extract the execution time  $e_t$  for each test  $t$  from the Surefire XML reports.

### 3.1.3. Model Training and Tuning

We use machine learning models to identify recurring patterns in source code features that correlate with performance testing coverage. By employing classification algorithms, we aim to predict the suitability of a method for performance testing based on these features. RICCARDO ▶ *We draw inspiration from the study conducted by Laaber [30], where the authors applied machine learning techniques to classify stable and unstable Go benchmarks. Their study analyzed 230 open-source Go projects, encompassing a total of 4,461 unique benchmarks. Similar to our work, they faced the challenge of working with an unbalanced dataset.* ◀ The process applied in this paper involves several steps: training the classifiers with a dataset that includes source code features of methods and coverage information as labels to indicate whether a method is performance tested. We also explore alternative class rebalancing techniques, utilize feature selection methods, and perform hyperparameter tuning to optimize the prediction accuracy of our models. Specifically, our process involves three progressively improving steps, as detailed below by referring to Fig. 1. First, we train the machine learning models using their default settings (RQ<sub>3.A</sub>). Next, we evaluate the impact of dataset preprocessing techniques, such as class rebalancing and feature selection, both in isolation and in combination (RQ<sub>3.B</sub>). Finally, we apply hyperparameter tuning to the top five model configurations, identified on the basis of prediction accuracy, considering the algorithm, class rebalancing technique, and feature selection technique (RQ<sub>3.C</sub>).

▷ 3A – *Model Training*: First, we train the machine learning models using their default parameters, as provided by the *scikit-learn* library<sup>5</sup>, without applying any dataset preprocessing. The dataset includes source code features for each method and coverage labels indicating whether a method is benchmarked. We evaluate the models using accuracy, precision, recall, F1 score, and balanced accuracy, as shown in Table 4 and discussed in Section 3.2.3. Due to the large amount of data (177,696 methods) and the severe class imbalance with some projects showing significant differences between benchmarked methods, we decided not to employ the k-fold cross-validation.

---

<sup>5</sup><https://scikit-learn.org/>

Table 3: Employed ML models

Model	Acronym	Algorithm
Random Forest	RF	Ensemble
ADA Boosting	ADA	Ensemble - Boosting
Multi-Layer Perceptron	MLP	Neural Network
Decision Tree	DT	Decision
k-Nearest Neighbor	kNN	Nearest Neighbor
Logistic Regression	LR	Linear Model
Gradient Boosting	GB	Ensemble - Boosting
Hist Gradient Boosting	HGB	Ensemble - Boosting
Linear Discriminant Analysis	LDA	Discriminant Analysis
Complement Naive Bayesian	CNB	Probabilistic - Bayesian
Naive Bayesian	NB	Probabilistic - Bayesian

To ensure fair experimentation, we apply an 80-20 test split to every experiment. For each model training session, we randomly remove 20% of the dataset for testing, while leaving the remaining 80% for training. Our study considers 11 well-known machine learning models that are listed in Table 3. These models have been selected on the basis of their previous usage in software performance research area [60].

▷ 3B – *Class Rebalancing*: Data obtained from the wrangling phase are used to assess the prediction capability of the considered models. Our dataset consists of 177,696 methods from 28 Java projects. Among these, 10,975 methods are benchmarked, while 166,721 methods are not. This represents a severe class imbalance, with benchmarked methods constituting only 6.58% of the total. To address this imbalance, we investigate two rebalancing techniques: oversampling and undersampling. For oversampling, we use the Synthetic Minority Over-sampling Technique (SMOTE) to increase the instances of the minority class [47]. SMOTE is a popular technique used during the training phase of machine learning models to handle imbalanced data. It generates new synthetic instances for the minority class by examining its nearest neighbors. For undersampling, we apply Random Undersampling (RUS) to reduce the instances of the majority class [48]. RUS is a straightforward method that reduces the amount of data in the majority class to balance the dataset. This approach identifies the most common class and reduces its size to match the minority class one, thus resulting in a balanced

dataset with fewer overall instances.

▷ **3B** – *Feature Selection*: This process includes removing co-linear and multi-co-linear features to avoid redundancy and improve model interpretability. We use Recursive Feature Elimination with Cross-Validation (RFECV) [49] and Auto Spearman [50] to identify the most relevant features for the prediction task. RFECV is a technique used during the feature selection phase of model training. It works by recursively considering progressively smaller sets of features. Starting with an ordered set of features ranked by importance, the less relevant features are pruned iteratively. This process continues until an optimal set of features is identified. By incorporating cross-validation, RFECV automatically finds the best set of features, thus ensuring that the selected features enhance model performance and generalization. Auto Spearman is another feature selection method that ranks features based on their Spearman correlation with the target variable. This approach helps in identifying features that have a strong relationship with the target, further refining the set of features used for prediction.

The application of feature selection techniques resulted in the selection of 22 features using RFECV and 38 features using Autospearman, out of the 44 total source code features.

▷ **3B** – *Configuration Selection*: We consider various configurations by combining different algorithms and dataset preprocessing techniques. Each configuration is denoted as a tuple  $(a, cr, fs)$ , where  $a$  represents one of the 11 algorithms in our study,  $cr$  indicates a class rebalancing method, and  $fs$  denotes a feature selection technique. Specifically,  $cr$  can be: no class rebalancing (original dataset), SMOTE, or RUS. Similarly,  $fs$  can be: no feature selection (original dataset), RFECV, or AutoSpearman. For each configuration, we train the model on the preprocessed training set and test it on the original (non-preprocessed) testing set. Our analysis involves training and testing 99 different configurations<sup>6</sup>, which we subsequently rank based on their F1 scores. We chose the F1 score because it balances the trade-off between false positives and false negatives, making it particularly useful in our context where both types of errors are significant.

▷ **3C** – *Hyperparameter Tuning*: We select the top-5 configurations as

---

<sup>6</sup>The 99 configurations are derived from applying 11 models with 3 sampling methods and 3 feature selection techniques.

produced by the previous phase and apply hyperparameter tuning to them. Hyperparameter tuning involves adjusting the parameters of the machine learning models that are set before the learning process begins and are not learned from the data. We use `RandomizedSearchCV`<sup>7</sup> from the *scikit-learn* library to explore a range of hyperparameter values and find the best combination for the used models. In particular, `RandomizedSearchCV` helps identify the best set of hyperparameters that maximize the model performance. For tuning, we employ the F1 score as the objective metric. The results of the tuning process provided the final top-5 performing configurations for predicting whether a method is benchmarked.

### 3.2. Employed metrics

Our investigation utilized various metrics, each pertinent to a specific research question as shown in Table 4, and detailed below.

#### 3.2.1. Coverage metrics ( $RQ_1$ )

We rely on method-level coverage to assess the coverage of performance (or functional) testing suites. Specifically, we employ four metrics: *code coverage*  $C_T$ , *direct code coverage*  $\hat{C}_T$ , *overlap ratio*  $OR_T$ , and *scope*  $S_T$ .

- *Code coverage* ( $C_T$ ) denotes the percentage of project methods covered—either directly or indirectly—by at least one test of the testing suite. Formally, a project method  $m \in \mathcal{M}$  is considered as *covered* by a testing suite  $T$  if there exist at least one test  $t \in T$  such that  $m \in M_t$ , where  $M_t$  denotes the set of methods invoked within the execution of  $t$ .
- *Direct code coverage* ( $\hat{C}_T$ ) is defined as the percentage of project methods that are *directly covered* by at least one test of the testing suite. A project method  $m$  is considered as *directly covered* by a testing suite  $T$  if there exist at least one test  $t \in T$  such that  $m \in \hat{M}_t$ , where  $\hat{M}_t$  denotes the set of methods directly invoked by  $t$ .
- *Overlap ratio* ( $OR_T$ ) measures the degree of redundancy within a testing suite  $T$ , by quantifying the extent of method coverage overlap across

---

<sup>7</sup>[https://scikit-learn.org/stable/modules/generated/sklearn.model\\_selection.RandomizedSearchCV.html](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.RandomizedSearchCV.html)

Table 4: Employed metrics

	<b>Employed metrics</b>		
<b>RQ<sub>1</sub></b>	<i>Coverage</i>	$C_T = \frac{ \bigcup_{t \in T} M_t }{ \mathcal{M} }$	$\hat{C}_T = \frac{ \bigcup_{t \in T} \hat{M}_t }{ \mathcal{M} }$
	<i>Overlap Ratio</i>	$OR_T = \frac{ \bigcup_{i,j \in T, i \neq j} (M_i \cap M_j) }{ \bigcup_{k \in T} M_k }$	
	<i>Scope</i>	$S_T = \frac{\sum_{t \in T}  M_t }{ T }$	
<b>RQ<sub>2</sub></b>	<i>Total Execution Time</i>	$TET_T = \sum_{t \in T} e_t$	
	<i>Average Execution Time</i>	$AET_T = \frac{\sum_{t \in T} e_t}{ T }$	
<b>RQ<sub>3</sub></b>	<i>Precision</i>	$\frac{TP}{TP+FP}$	
	<i>Recall</i>	$\frac{TP}{TP+FN}$	
	<i>F1 Score</i>	$\frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$	
	<i>Balanced Accuracy</i>	$\frac{\text{Recall of Class 1} + \text{Recall of Class 2}}{2}$	

different tests. Table 4 provides a formal definition of this metric, where  $i$  and  $j$  denote two distinct tests that belong to  $T$ ,  $M_i$  represents the set of methods covered by test  $i$ , and  $M_j$  represents the set of methods covered by test  $j$ . The numerator in  $OR_T$  denotes the number of methods covered in more than one test, while the denominator denotes the number of methods covered by the testing suite.  $OR_T$  values range from 0 to 1, where 0 indicates no overlap, i.e., each test cover distinct methods, and 1 indicates high test redundancy, i.e., all the methods are covered by more than one test.

- *Scope* ( $S_T$ ) measures the average number of methods that an individual test covers within a testing suite. We use this metric because previous work has demonstrated that high coverage of tests (i.e., scope) tends to have a positive impact on the capabilities of uncovering performance issues [11].

### 3.2.2. Time cost metrics (RQ<sub>2</sub>)

For each testing suite  $T$ , we consider two time cost metrics: *total execution time* ( $TET_T$ ) and *average execution time* ( $AET_T$ ). The *total execution time* ( $TET_T$ ) is the cumulative sum of the execution times for all tests in the suite, i.e., the time required to run the entire testing suite. The *average execution time* ( $AET_T$ ) is calculated as the mean execution time of the individual tests within the suite.

### 3.2.3. Model performance metrics (RQ<sub>3</sub>)

To evaluate the prediction accuracy of our machine learning models, we employ well-established metrics for binary classification: *precision*, *recall*, *F1 score*, and *balanced accuracy*. We specifically choose *balanced accuracy* over traditional *accuracy* because the latter can be misleading in the presence of class imbalance [68].

## 4. Results and Discussion

In this section, we present and discuss the results of our analysis. As illustrated in Table 5, we analyzed 2,190 JMH benchmarks from 28 software projects, and 256,233 JUnit tests from 19 software projects. For RQ1 and RQ2, we analyze performance tests on all 28 projects (both for coverage and time cost). However, 9 Projects were excluded from the comparison of performance and functional testing suites due to technical issues encountered during the JUnit data collection (see Section 3.1 for details).

### 4.1. RQ<sub>1</sub>: To what extent do performance tests cover the source code of software systems?

To answer RQ<sub>1</sub>, we analyze the coverage of JMH benchmarks within our dataset of 28 software projects. We centre our analysis around two key metrics: (i) *Benchmark Coverage* ( $C_T$ ), which measures the extent of method coverage by JMH benchmarks, and it includes *direct coverage* ( $\hat{C}_T$ ); (ii) *Overlap Ratio* ( $OR_T$ ), which assesses the degree of redundancy in method coverage across different benchmarks of the same project. Besides this, we analyze the coverage of JUnit tests and subsequently compare it with the coverage of JMH benchmarks within a subset of 19 software projects, as mentioned above.

*Benchmark Coverage.* Our analysis revealed that the coverage by JMH benchmarks in software projects is relatively low compared to the total number of methods:  $C_T$  averaged 8.8% with 2.1% of methods directly covered ( $\hat{C}_T$ ) across all 28 projects.

Figure 3 shows the distribution of benchmark coverage ( $C_T$ ) across the examined projects. The y-axis represents the percentage of methods benchmarked in each project, while the x-axis enumerates the projects. The project with the highest coverage is **panda**, where 48.82% of methods are covered. On the other hand, the **jooby** project has the lowest coverage (0.27%). The analysis indicates diverse levels of coverage among the systems with a standard deviation of 9.2%. By excluding the **panda** project, considered as an outlier, from our analysis, the benchmark coverage’s standard deviation was 4.9%.

Table 5: Java systems overview

Project	GitHub Stars	Methods (Total)	Benchmarks (Total)	Unit Tests (Total)	Domain
arrow	12600	3789	34	869	Analytics Tools
byte-buddy	5800	5046	39	6357	Code Generation
cantaloupe	259	3475	103	3063	Computer Graphics
client_java	2100	386	33	217	JVM Tools
commons-bcel	223	3132	3	137	JVM Tools
crate	3800	24155	39	-	Database Systems
eclipse-collections	2300	15219	515	-	Programming Utility
fastjson	25500	17524	4	4979	Parsing Library
feign	9100	979	8	913	Web Development
HdrHistogram	2100	780	12	147	Analytics Tools
imglib2	278	3432	25	635	Computer Graphics
iri	1200	1554	3	398	Data Structures
jdbi	1800	2379	76	1428	Database Systems
jetty.project	3700	18060	48	-	Web Development
jgrapht	2400	3782	51	2416	Programming Utility
jooby	1600	3331	3	485	Web Development
kafka	26000	16157	27	-	Data Streaming
logbook	1600	811	20	564	Web Development
netty	31900	16615	221	-	Network Applications
objenesis	568	207	13	45	Programming Utility
panda	247	1479	4	61	Programming Utility
protostuff	2000	3312	16	-	Programming Utility
r2dbc-h2	191	291	8	259	Database Systems
rdf4j	331	14041	14	-	Database Systems
RxJava	47300	8282	217	-	Programming Utility
SquidLib	439	5663	236	73	Computer Graphics
tinkerpop	1800	7753	57	-	Database Systems
vert.x	13800	5685	41	4095	JVM Tools

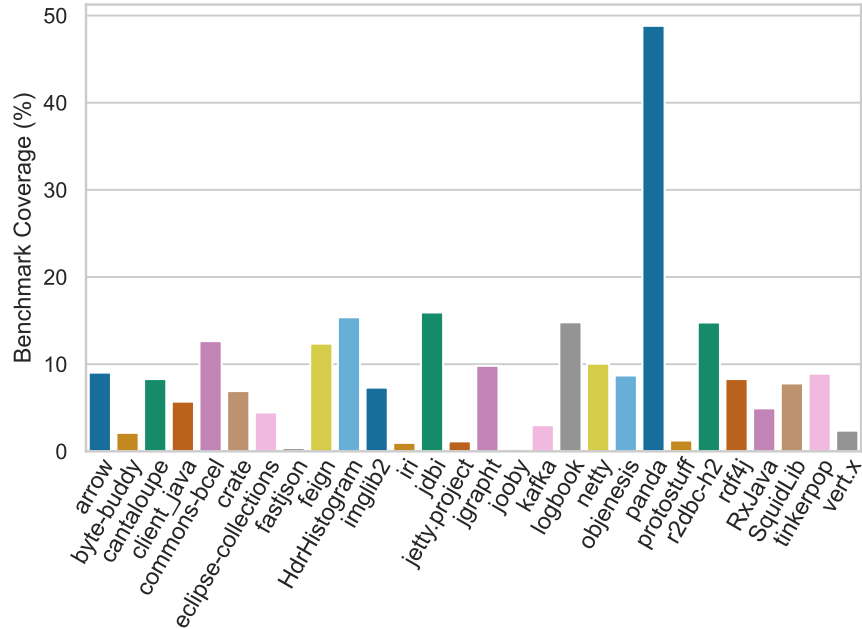


Figure 3: Coverage ( $C_T$ ) of benchmarks across projects.

This variation can be attributed to differences in the size of the projects (in terms of their number of methods, benchmarks and unit tests), to the particular development and testing practices adopted, and to the specific performance requirements and constraints. It is evident from the figure that **panda** is an outlier. Perhaps, in **panda**, we can associate the high coverage of benchmarks with the nature of the project, which is a lightweight and extensible programming language toolkit. Due to required efficiency and scalability, developers may focus extensively on microbenchmarking.

A detailed look at *direct benchmark coverage*, depicted in Figure 4, reveals a similar trend of varied coverage. As a first consideration, the average direct benchmark coverage across many projects is around 2%, substantially lower than the overall benchmark coverage. The figure also shows a sensibly lower variation in the direct coverage compared to the overall coverage. The standard deviation for direct benchmark coverage is 1.7%, which also indicates a lower spread than the one in overall benchmark coverage. The observed difference between direct and indirect coverage may be explained by developers' tendencies to target high-level methods during performance testing. Such high-level methods often result in a large number of indirect

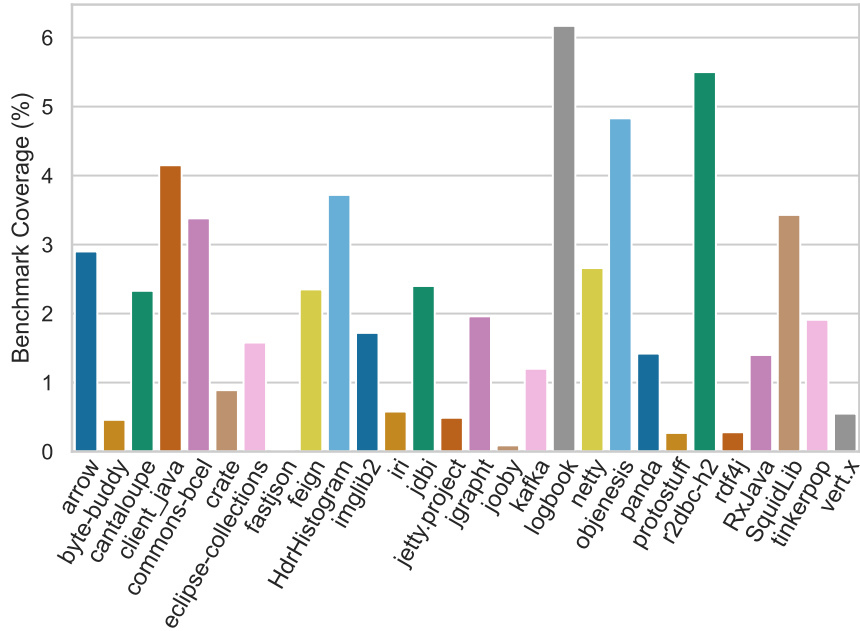


Figure 4: Direct coverage ( $\hat{C}_T$ ) of benchmarks across projects.

invocations, which might broaden the overall benchmark coverage.

*Overlap Ratio.* Figure 5 illustrates the degree of overlap in method coverage, which represents the redundancy in method coverage across different benchmarks in the project.

The overlap ratio in the considered projects ranges from a minimum of 0.09 in `commons-bcel` to a maximum of 1 in `byte-buddy`. In general, we can observe relatively high overlaps, with an average of 74% methods covered by more than one benchmark. While this average suggests a high degree of redundancy within performance testing suites, this may also reflect the attempts of developers to test a method under various workloads [69]. Nonetheless, our analysis highlights room for improving the efficiency of performance testing by reducing unnecessary overlap.

*Comparison of JMH Benchmarks and JUnit Tests Coverage.* Figure 6 illustrates a comparison between the coverage achieved by performance testing and the functional testing one. Clearly, the percentage of methods covered by JUnit tests tends to be much higher than the one of JMH benchmarks,

thus indicating a broader coverage for functional testing in the projects under study. On average, the coverage achieved by a performance testing suite is 4 times less than a functional testing suite one (10.4% *versus* 41.3%). The JUnit tests coverage is significantly higher in many projects (like `arrow`, `cantaloupe`, `fastjson`, `iri`, `jooby`, `jgrapht`, and `vert.x`) as compared to the JMH benchmarks coverage. However, it is also interesting to note an exception, i.e., the `SquidLib` project, where not only the JUnit Test coverage is relatively low, but also the JMH benchmark coverage exceeds the JUnit tests one. Upon closer examination of the project code, we attributed this anomaly to the performance-driven focus of the project. Specifically, `SquidLib` is a Java library crafted to serve as a comprehensive toolkit for the development of various gaming applications. This orientation likely leads to a higher priority being placed on performance testing over unit testing, thereby resulting in a more extensive benchmark coverage.

Similarly, Figure 7 provides a comparison of direct coverage between JMH benchmarks and JUnit tests across the same set of projects. The trend is similar to the previous one. Few exceptions appear also here, such as (again)

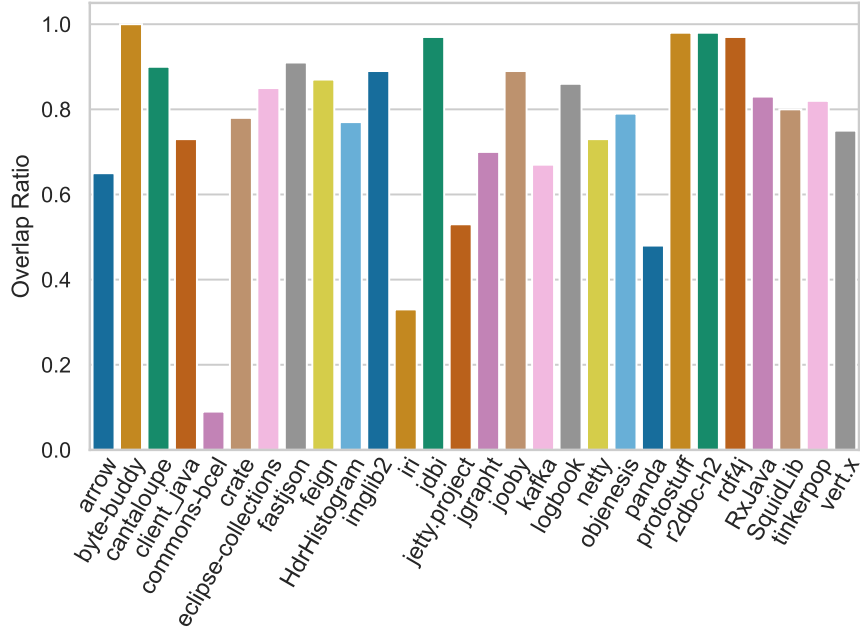


Figure 5: Overlap Ratio ( $OR_T$ ) of benchmarks.

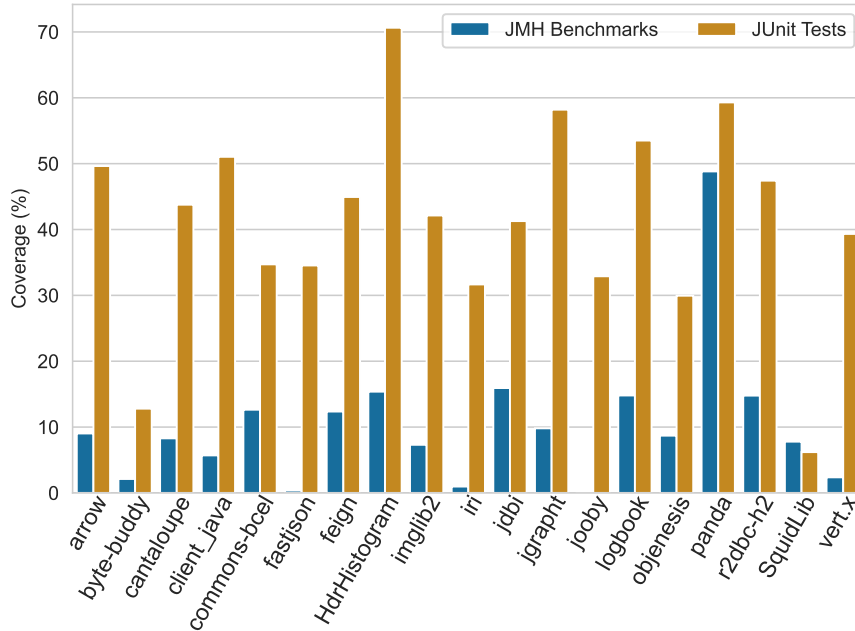


Figure 6: Comparison of coverage ( $C_T$ ) of JMH Benchmarks and JUnit Tests across projects.

SquidLib, where direct benchmark coverage exceeds the direct coverage by unit tests. We also note substantial differences across projects in terms of the direct coverage of methods by both JMH benchmarks and JUnit tests. For instance, in the `r2dbc-h2` project, the direct coverage of JUnit tests reaches 30%, whereas the one of JMH benchmarks is nearly 5%. This observation emphasizes the idea that, while functional tests point to broad coverage to ensure overall correctness, performance benchmarks often target specific, performance-sensitive parts of the code.

*Overlap Ratio Comparison.* Our analysis reveals that JMH benchmarks exhibit a more significant overlap in coverage compared to JUnit tests one. In particular, we found that JMH suites show higher overlap in coverage compared to that of the JUnit suites in 68% of the projects. As shown in Figure 8, `byte-buddy`, `r2dbc-h2`, and `jdbi` exhibit near-complete overlap for JMH benchmarks, indicating extensive redundancy in performance testing. On the other hand, `commons-bcel` stands out with a much lower overlap ratio for JMH benchmarks (0.09), yet a higher overlap for JUnit tests (0.66), indi-

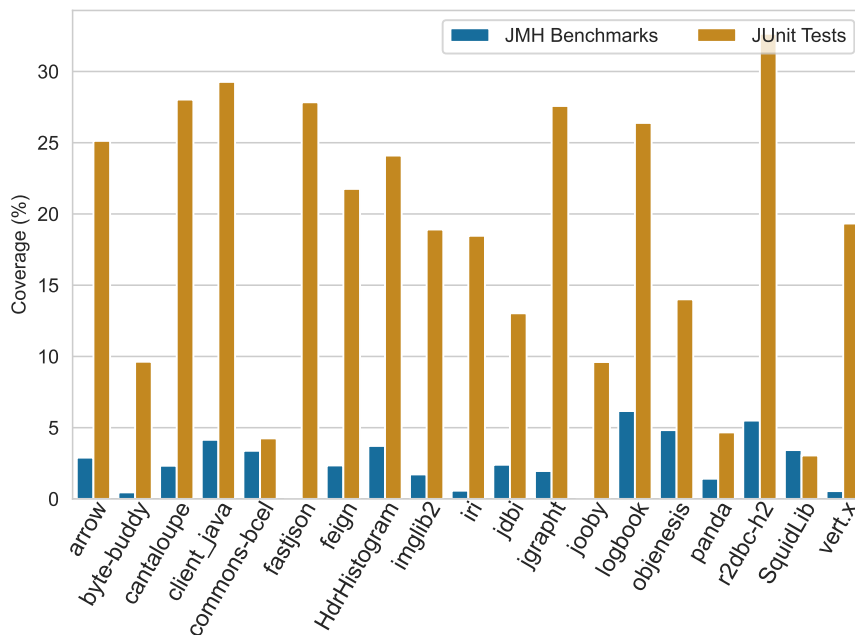


Figure 7: Comparison of direct coverage ( $\hat{C}_T$ ) of JMH Benchmarks and JUnit Tests across projects.

ating a different testing approach. Meanwhile, `fastjson` has a high overlap in JMH benchmarks (0.91), but a much lower ratio in JUnit tests (0.15). These variations highlight the different testing priorities and strategies between performance and functional correctness across projects.

*Scope Comparison of JMH Benchmarks and JUnit Tests.* Figure 9 compares the scope ( $S_T$ ) of JMH benchmarks and the JUnit tests one. It is interesting to observe that, while the coverage of performance testing suites is generally lower than the functional testing suites one, the average coverage (i.e., scope) of individual benchmarks is notably larger. Specifically, JMH benchmarks show, on average, approximately three times larger scope than their JUnit counterparts. This significant difference in the scope is more evident in projects like `panda` and `commons-bcel`, likely due to the project-specific nature of JMH benchmarks, which emphasize testing performance under various configurations and workloads. These results, once again, suggest that JMH benchmarks tend to target high-level methods, which leads to a larger number of method invocations, whereas JUnit tests tend to target

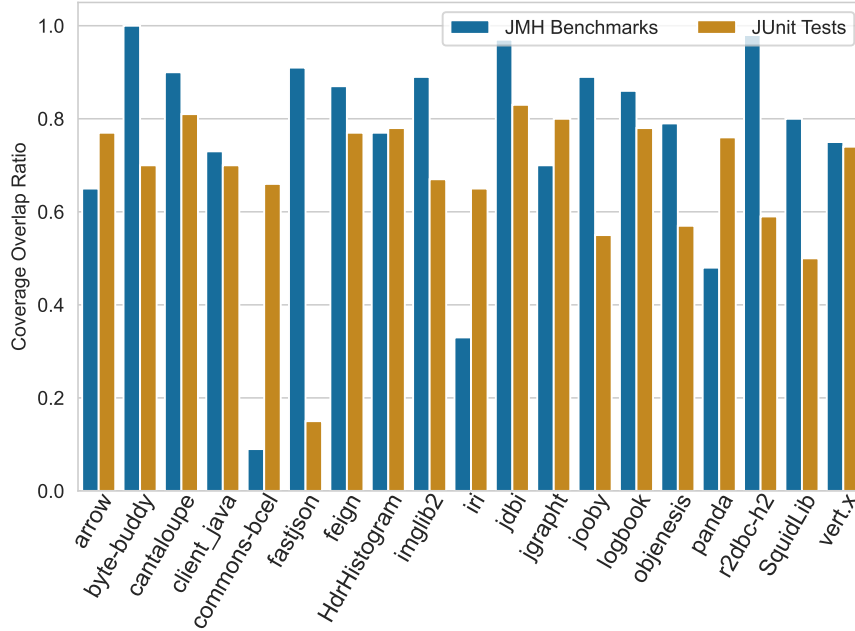


Figure 8: Comparison of Overlap Ratio coverage ( $\hat{OR}_T$ ) of JMh Benchmarks and JUnit Tests across projects.

lower-level methods. This finding is in line with our previous observation about the difference between direct and indirect coverage.

**Answer to RQ<sub>1</sub>:** The analyzed projects show that performance testing suites achieve a limited code coverage of 8.8% on average, which is four times lower than that of functional tests. Additionally, in 68% of the projects, performance testing suites exhibit higher overlap than functional testing suites, thus indicating larger redundancy in the tested methods.

#### 4.2. RQ<sub>2</sub>: What is the time cost of performance tests?

In this subsection, we present our findings on the execution time of JMh benchmarks and compare them to the functional testing ones.

*Total Execution Time.* Figure 10 displays the execution time of performance testing suites. On the y-axis, we report the total execution time ( $TET_T$ )

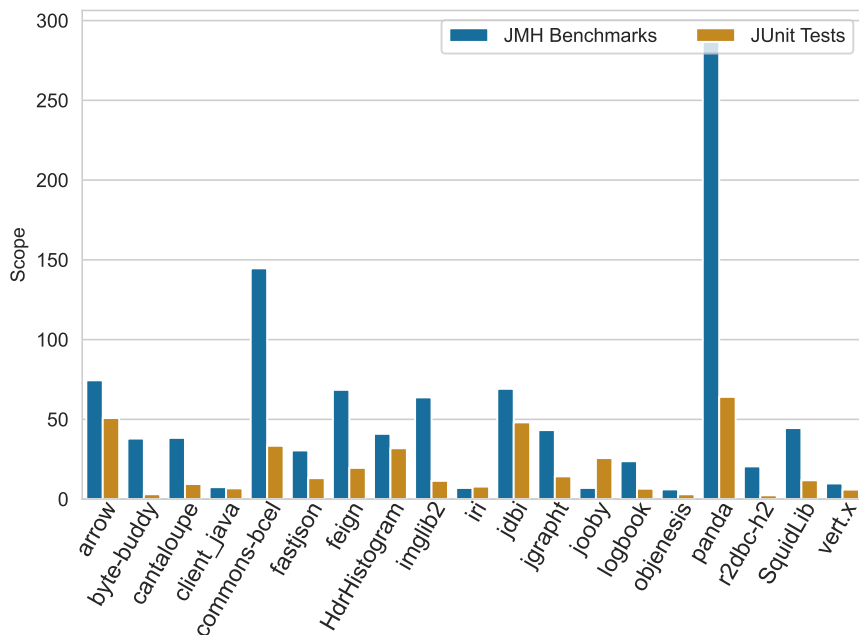


Figure 9: Scope ( $S_T$ ) of JMH Benchmarks vs. JUnit Tests.

in hours, using a logarithmic scale. The black dashed line depicts the overall average time, which is about 29.3 hours. The  $TET_T$  distribution varies significantly from one project to another, with some testing suites completing in just a few minutes, while others requiring over 100 hours of execution. The less time-consuming suite is the `commons-bcel` one, which required only 90 seconds to run its benchmarks. Along with `commons-bcel`, only other two projects kept the  $TET_T$  under 6 minutes (i.e.,  $10^{-1}$  hours), namely `panda` and `r2dbc-h2`. About the most time-consuming performance testing suites, few ones exceeded 10 hours. Interestingly, two projects (i.e., `squidLib` and `rxJava`) exceeded 10 hours but did not overcome the threshold of 100 hours, which was required by other three projects. In particular, `eclipse-collection` required 401 hours, thus representing the most time-intensive project for performance testing.

*Average Execution Time.* In Figure 11, the bar chart depicts the average execution time of benchmarks ( $AET_T$ ) for each performance testing suite. The dashed line shows the average  $AET_T$  across projects, which is about 23 minutes. A close examination of this chart confirms that the diver-

sity across projects observed for  $TET_T$  also holds for  $AET_T$ . In particular, `commons-bcel`, which was the least time consuming one in terms of  $TET_T$ , resulted in a relatively low  $AET_T$  of 30 seconds, and similarly `panda` and `r2dbc-h2` have maintained a low average execution time. Contrariwise, `eclipse-collections` (515 benchmarks and an  $AET_T$  of 2,805 seconds) and `netty` (221 benchmarks and an  $AET_T$  of 2,149 seconds) are very time-consuming projects, as equipped with a high amount of benchmarks.

The most time-consuming performance testing suite is the one of `kafka`, with an  $AET_T$  of about 7 hours per benchmark, followed by `eclipse-collections`. Interestingly, the performance testing suite of `kafka` consists of a limited number of benchmarks (i.e., 27), each of which is then notably time-consuming. We investigated the JMH reports to understand the reasons behind these high  $AET_T$  values, and we discovered that the likely reason is the large number of repetitions of `kafka` benchmarks under different input values. This allows the benchmark to measure performance across a variety of configurations. Each configuration results in an additional execution of the benchmark, which naturally increases the total execution time [69].

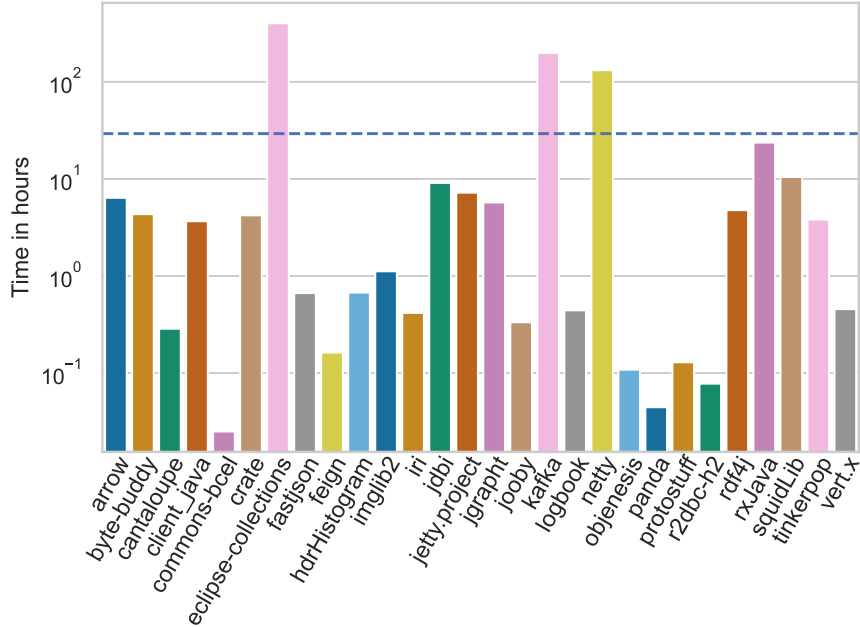


Figure 10: JMH total execution time ( $TET_T$ ).

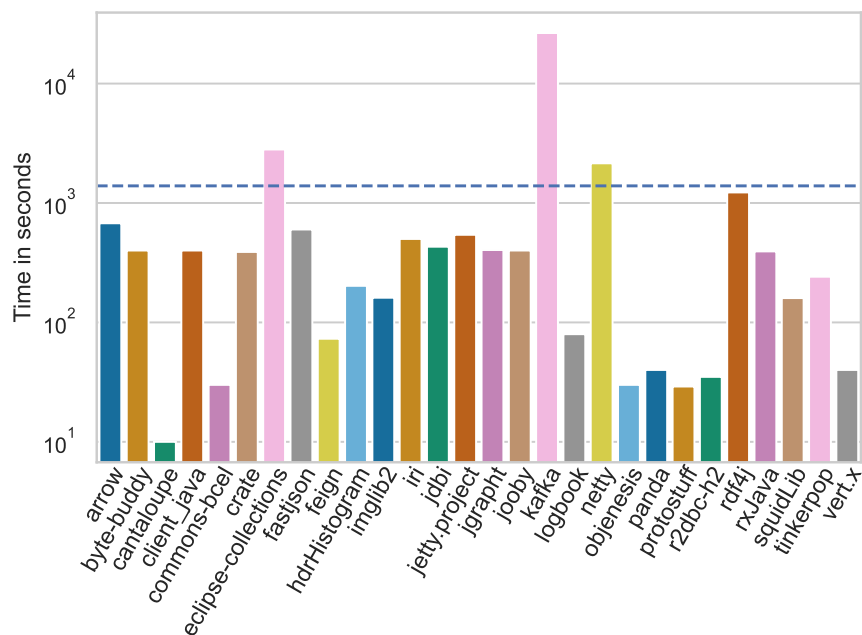


Figure 11: JMHAverage execution time ( $AET_T$ ).

*Total Execution Time Comparison.*: In Figure 12, we compare  $TET_T$  of JMHA and JUnit testing suites. As expected, the analysis revealed that, in general, performance testing is significantly more time-consuming than functional testing. For instance, by examining `SquidLib`, the most time-consuming project for performance testing, we observe that  $TET_T$  for the JMHA suite is exponentially higher than the JUnit suite one (37,735 seconds *versus* 18.5 seconds). Interestingly, the least time-consuming project in terms of functional testing is `vert.x` that requires 821 seconds, whereas  $TET_T$  for its JMHA suite is 1,640 seconds, i.e., more than twice the JUnit test suite execution. We can notice comparable results if we analyze less time-consuming projects. For instance, `commons-bcel` requires 90 seconds for executing the whole JMHA suite and reports a  $TET_T$  of 15.27 seconds for the JUnit suite. Another interesting case is `objenesis`, which requires half a second for executing the JUnit testing suite, and 390 seconds for executing the JMHA suite. Performance testing suites have, on average, an execution time cost 62 times higher than functional testing suites one.

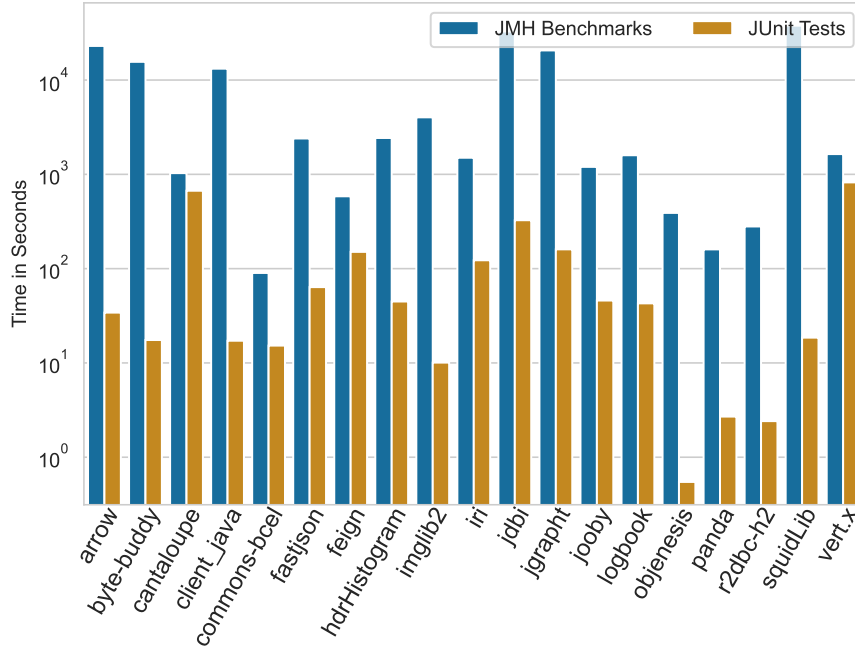


Figure 12: Total execution time ( $TET_T$ ) comparison.

*Average Execution Time Comparison.*: There is a significant difference in test-level execution time between JMH and JUnit tests across all projects. We found that JMH benchmarks exhibit an average  $AET_T$  of approximately 245.9 seconds. On the other side, the average  $AET_T$  of JUnit tests does not exceed 0.098 seconds. On average, an individual JMH benchmark demands about 2,507 times the execution time required by a JUnit test. We report the complete results related to  $AET_T$  of JUnit tests in our replication package [19].

**Answer to RQ<sub>2</sub>:** In the analyzed projects, performance testing suites have, on average, a time cost 62 times higher than that of functional testing suites. Specifically, the average execution time for a performance testing suite is 29.3 hours, while individual performance tests last 23 minutes on average. These results highlight the significant execution time cost of performance testing, particularly when compared to the much shorter execution times of functional testing.

4.3. *RQ<sub>3</sub>: Are there any specific characteristics in the source code that can guide performance testing?*

To address RQ<sub>3</sub>, we examine the relationship between source code features and the likelihood of being covered by performance tests. Our goal is to search for peculiarities in the methods selected for performance testing. To achieve this, we investigate the possibility of employing machine learning models to predict whether a method should be performance-tested based on its source code features.

4.3.1. *RQ<sub>3,A</sub>: Can we predict whether a method should be performance tested using machine learning models with default settings?*

We first analyze the prediction accuracy of 11 classification algorithms (discussed in section 3.1.3) using their default settings.

Figure 13 presents the results for the employed machine learning models in terms of various performance metrics, including precision, recall, F1 score, and balanced accuracy, as detailed in Table 4. Each subplot corresponds to a different performance metric used to evaluate prediction performance. The results indicate that it is challenging to predict with high accuracy which methods should be benchmarked based on source code features. The F1 score, which is particularly meaningful in this context, shows that the best models barely reached 0.4%, with none of the tested models exceeding 0.5%. Furthermore, only the Decision Tree, k-Nearest Neighbors, and Random Forest models achieved reasonable performance levels. In contrast, the Naïve Bayesian (NB) model performed the worst, remaining below 0.2%.

As shown in Figure 13, the Naïve Bayesian model demonstrates a high recall but very low precision. This indicates a higher rate of false positives, meaning that while the NB model is effective at identifying methods that should be benchmarked (high recall), it also incorrectly flags many unsuitable methods as suitable (low precision). This leads to a higher rate of false positives, thus making the model less reliable when it is crucial to ensure that most of the predicted positives are true positives.

On the other hand, the Gradient Boosting (GB) and Histogram Gradient Boosting (HGB) models perform apparently well in terms of precision. This score is due to the small amount of true and false positive predictions. However, if we consider their recall, then we notice values close to 0, thus confirming that most of the predictions were true negative. This fact is proved also by the F1 score values, close and lower than 0.1%.

F1 score balances precision and recall, thus providing a comprehensive metric to evaluate models. Consequently, the Decision Tree (DT), k-Nearest Neighbors (kNN), and Random Forest (RF) models, with their higher F1 scores, appear the most effective ones for predicting whether a method should be benchmarked based on static features of the source code. However, it is important to note that the absolute F1 scores of these models are quite low, thus implying significant room for improvement in their prediction capabilities.

Since we are dealing with a relevantly unbalanced dataset, as reported by the coverage analysis, we have also employed the Balanced Accuracy metric. The results of this latter metric show that, in all the examined cases, the models can correctly identify the true negatives (i.e., methods that are correctly classified as not benchmarkable). This result is due to the evident unbalanced situation. In other words, the models easily identify the non-benchmarkable methods because, during the training phase, most of the methods were not covered.

*4.3.2. RQ<sub>3.B</sub>: How does the prediction accuracy change when employing feature selection and class-rebalancing techniques?*

To answer RQ<sub>3.B</sub>, we investigate whether and to what extent pre-processing steps impact the model prediction performance. To do this, we apply a combination of two commonly used pre-processing steps: (1) class-rebalancing using the Synthetic Minority Over-sampling Technique (SMOTE) [47] and Random Under-Sampling (RUS) [48], and (2) removal of co-linear and multi-co-linear features using AutoSpearman [50] and Recursive Feature Elimination with Cross-Validation (RFECV) [49].

Figure 14 illustrates the prediction performance metrics for all studied algorithms across various combinations of these pre-processing techniques. We assess nine distinct configurations: (1) a baseline without feature selection and class rebalancing, (2) RFECV without rebalancing, (3) AutoSpearman without rebalancing, (4) no feature selection with SMOTE class rebalancing, (5) no feature selection with RUS, (6) RFECV combined with SMOTE, (7) AutoSpearman combined with SMOTE, (8) RFECV combined with RUS, and (9) AutoSpearman combined with RUS. This comprehensive approach results in 99 distinct configurations (11 models  $\times$  9 pre-processing combinations), thus allowing us to thoroughly evaluate the impact of various pre-processing techniques on model performance.

Each subplot in Figure 14 represents the F1 scores of a specific machine

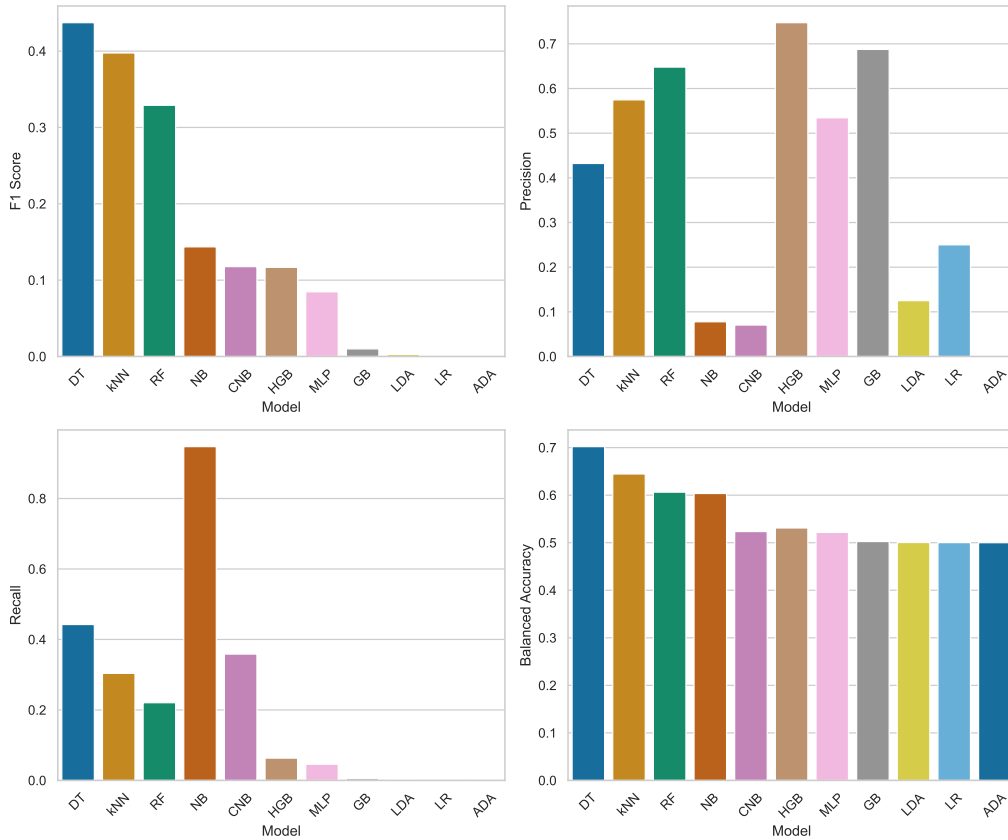


Figure 13: Predictive Performance of Static Source Code Features

learning model across different combinations of feature selection and sampling techniques. Significant variations in model performance across different pre-processing configurations can indeed be observed. First of all, the effect of removing co-linear and multi-co-linear features using RFECV and AutoSpearman is a notable variation across models. For example, Random Forest (RF) always shows a slight improvement in F1 score when RFECV is applied (from 0.329 to 0.340 without sampling and from 0.296 to 0.298 with RUS). It has to be noticed that, when combined with SMOTE, RF's performance increases to reach an F1 score of approximately 0.47. This indicates that RF benefits from both feature selection and class rebalancing in general (specifically RFECV and SMOTE).

The subplot for k-Nearest Neighbors (kNN) in Figure 14 shows that the

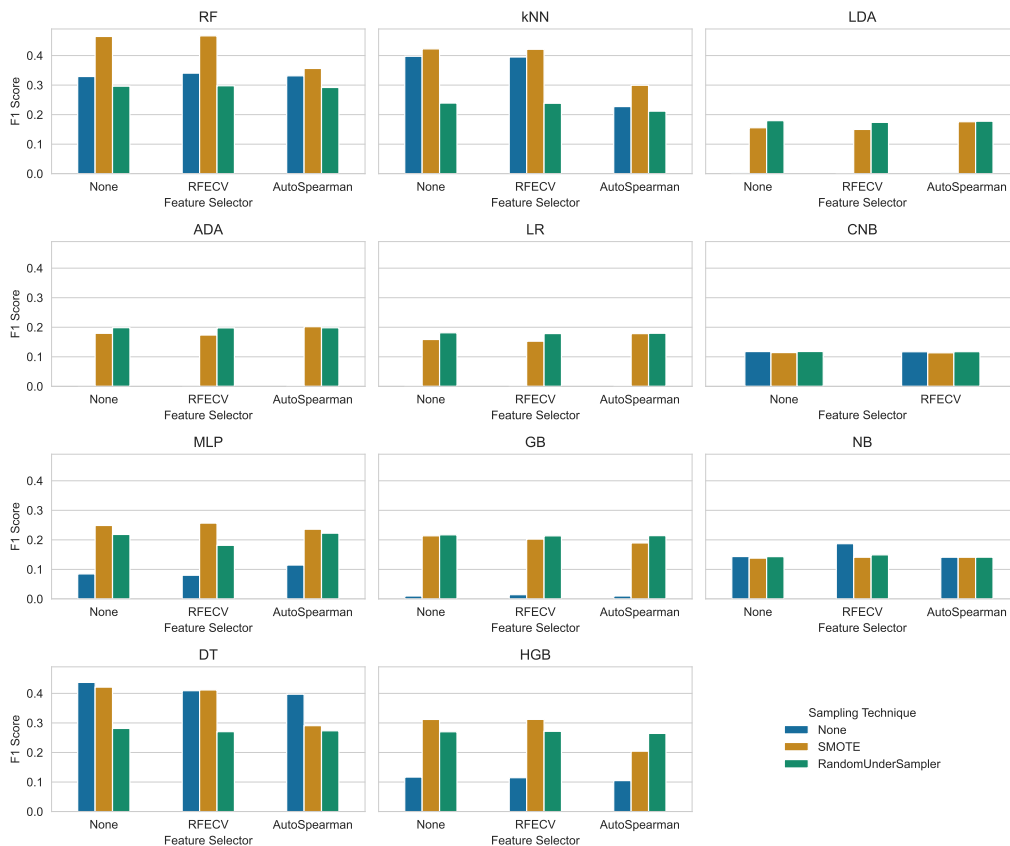


Figure 14: F1 Scores for All Models with Sampling Methods and Feature Selection Techniques

baseline F1 score without any pre-processing is better than the one obtained by using AutoSpearman for feature removal. In addition, the usage of RFECV does not provide any significant improvement. The combination of RFECV and SMOTE results in an F1 score of 0.421, which, while higher than the baseline, does not overcome the performance of SMOTE alone, which achieves the highest F1 score for kNN at 0.422.

For Decision Trees (DT), the highest baseline F1 score is 0.437. The application of RFECV alone slightly reduces the score to 0.409, while AutoSpearman lowers it further to 0.397. Combining RFECV with SMOTE performs relatively better, but the baseline is still more effective. However, for DT, SMOTE also occurs to be more effective than RUS with all the

feature pre-processing techniques.

In general, models such as Logistic Regression (LR), Gradient Boosting (GB), and Linear Discriminant Analysis (LDA) show improvements with pre-processing. The baseline F1 score for LR is nearly zero, and the application of RUS alone provides the highest F1 score. This shows the effect of class rebalancing on the performance improvement of this model.

For Naive Bayes (NB), the application of RFECV alone provides the highest F1 score of 0.187, whereas Complement Naive Bayes (CNB) exhibits relatively stable performance across different pre-processing configurations. CNB shows minimal variation in performance, with F1 scores remaining around 0.113 to 0.118.

Therefore, in answering RQ<sub>3.B</sub>, our analysis shows that the impact of removing co-linear and multi-co-linear features, as well as applying class-rebalancing techniques, significantly varies across different machine learning models. Feature selection methods like RFECV and AutoSpearman exhibit model-specific effects, with some algorithms benefiting more than others. Class-rebalancing techniques, particularly SMOTE, generally show improvement in model performance for algorithms like Random Forest and k-Nearest Neighbors. However, Decision Trees perform better with the baseline configuration, while models like Logistic Regression and Naive Bayes show only incremental improvements. Thus, the combination of feature selection and class-rebalancing does not always yield additive benefits.

It is important to note that, despite these pre-processing steps, the overall prediction performance based on the static source code features remains relatively low across all models and configurations. The highest F1 score achieved is only 0.467, which indicates that there is still significant room for improvement in predicting the likelihood of a method being tested for performance. This suggests that while pre-processing steps can enhance model performance to some extent, they are not sufficient to overcome the challenges in this prediction task.

#### *4.3.3. RQ<sub>3.C</sub>: How does the prediction accuracy change after tuning the hyperparameters of machine learning models?*

Based on findings from RQ<sub>3.A</sub> and RQ<sub>3.B</sub>, in this subsection, we answer RQ<sub>3.C</sub> by selecting the top 5 configurations with the highest F1 scores from the previous analysis (see Table 6). We then apply hyperparameter tuning to the respective models in these top 5 configurations to evaluate the effect on prediction performance (see Table 7).

By comparing Table 6 and Table 7, it is evident that hyperparameter tuning resulted in moderate improvements in model performance. The best F1 score increased slightly from 0.47 to 0.48 for the configuration of Random Forest with SMOTE and RFECV. Also, Random Forest maintained its position as the best-performing model both before and after hyperparameter tuning, which shows the relatively robust configuration for predicting the likelihood of the method to be covered by performance testing. Interestingly, while the top two configurations remain consistent (RF with SMOTE), we observe changes in the subsequent rankings. The k-Nearest Neighbors model with SMOTE and RFECV shows some improvement, with its F1 score increasing from 0.42 to 0.46 after tuning. Similarly, hyperparameter tuning affected precision and recall differently across models. For instance, the baseline configuration of Decision Tree model (no sampling and no feature selection) had an improvement in recall (from 0.44 to 0.50) but a decrease in precision (from 0.43 to 0.37). Most models showed a slight increase in Balanced Accuracy after hyperparameter tuning with the exception of kNN where it decreased from 0.80 to 0.76.

Therefore, the application of hyperparameter tuning has not considerably improved the overall prediction performance, with the highest F1 score reaching only 0.48. We observe that while hyperparameter optimization improves the performance of some models with certain configurations compared to others, it still lacks promising results for this specific task of identifying methods that should be benchmarked based on static code features. This also highlights that predicting which methods should be benchmarked based on static code features remains a challenging task.

Table 6: Top 5 Configurations Based on F1 Score without Hyperparameter Tuning

Model	Sampling Technique	Feature Selector	Precision	Recall	Balanced Accuracy	F1 Score
RF	SMOTE	RFECV	0.56	0.40	0.69	0.47
RF	SMOTE	None	0.56	0.40	0.69	0.46
DT	None	None	0.43	0.44	0.70	0.44
kNN	SMOTE	RFECV	0.30	0.71	0.80	0.42
DT	SMOTE	None	0.37	0.49	0.71	0.42

Table 7: Top 5 Configurations Based on F1 Score with Hyperparameter Tuning

Model	Sampling Technique	Feature Selector	Precision	Recall	Balanced Accuracy	F1 Score
RF	SMOTE	RFECV	0.55	0.42	0.70	0.48
RF	SMOTE	None	0.55	0.40	0.69	0.46
kNN	SMOTE	RFECV	0.39	0.58	0.76	0.46
DT	None	None	0.37	0.50	0.70	0.43
DT	SMOTE	None	0.37	0.48	0.71	0.42

**Answer to RQ<sub>3</sub>:** In the analyzed projects, machine learning models exhibit limited accuracy in predicting whether a method is subject to performance testing, by achieving a top F1 score of 0.44 with default settings. Even after applying well-known practices, such as feature selection, class rebalancing, and hyper-parameter tuning, the best model only achieves a top F1 score of 0.48. These results do not reveal any clear relationship between recurrent source code patterns in methods and the likelihood of being performance tested, thus suggesting a lack of generalizable source code characteristics to guide performance testing.

## 5. Implications

This section discusses some implications of this study along with some directions for future work.

**For practitioners.** This study gives clear evidence that a significant portion of the codebase of many software systems lacks performance assessment. As software evolves, these code areas may become vulnerable to performance bugs that remain undetected until released. One potential reason for this oversight could be the limited availability of tools for performance test coverage. Indeed, we are unaware of any tool that measures performance test coverage as seamlessly as tools like JaCoCo does for functional testing. We believe that the introduction of such a tool could allow developers to more regularly assess the coverage of their performance testing suites, thus increasing their awareness of the code area that remains unmonitored for performance.

Our results also reveal that distinct performance tests often cover the same code components (i.e., high overlap), even though a significant portion of the codebase remains uncovered by performance test. Although developers might deliberately target the same code components with multiple performance tests to assess their behaviour under varying workloads, this highlights an opportunity to broaden test coverage without incurring additional time costs. We hypothesize that, by raising awareness about performance test coverage, developers might be more inclined to prioritize the creation of tests that target code components currently unassessed for performance. We encourage future work to make it easier for practitioners to be aware of the performance testing suites coverage.

**For researchers.** Prior work suggests that the high costs of test development and maintainability often hamper the adoption of performance testing in practice [12, 5, 46]. In response to this issue, researchers have introduced techniques capable of automatically transforming functional testing suites into performance tests [12]. In light of our results, we can formulate educated guesses regarding the potential benefits of these techniques, as well as the challenges that might originate from their adoption. For instance, our results suggest that, by utilizing automated performance test generation, there could be a significant improvement in terms of performance test coverage. Indeed, functional testing suites exhibit significantly higher coverage than that of performance tests (10.4% *vs* 41.3% on average in our dataset), and generated performance tests would inherit such high coverage. However, these benefits might come at a cost, particularly regarding execution time. The high coverage of functional test suites is typically a consequence of their extensive sizes, which may not be feasible for a performance testing suite. In fact, performance tests are typically more time-consuming than functional counterparts (on average 2,507 times more in our dataset). For instance, by using the configuration defined by Jangali *et al.* [12] and a typical number of five forks [54, 52], the time cost of an individual generated performance test would amount to about 400 seconds. For a medium-sized testing suite like `logbook`, which comprises 564 JUnit tests, this translates to a total time cost of roughly two and a half days. This is approximately 141 times longer than the actual `logbook` performance testing suite. Even when considering the smallest functional testing suite in our study, namely `objenesis`, this results in a 5-hour execution time, i.e., 55 times the one of the current performance testing suite. These findings highlight that automated generation alone might

not be sufficient to produce performance test suites that are practically usable, given that developers might be deterred by such a time-consuming test process. This underscores a significant challenge for the research community, i.e., automated performance test generation should take into consideration the associated time cost of the generated testing suite.

A potential research direction we aimed to trigger is to devise “smart” selection strategies that automatically identify which code components should undergo performance testing. However, our results indicate that developing a selection strategy aligned with developers’ perceptions of what should be tested can be challenging. Our findings highlight a lack of recurrent patterns in source code subject to performance testing, thus suggesting that performance testing should be driven primarily by domain-specific objectives rather than generalizable source code features across software systems. Consequently, researchers should consider other viable alternatives. One option could be to employ test selection strategies that aim to maximize code coverage while mitigating the time cost, for instance, by reducing the number of redundant performance tests covering the same code component [14, 15]. Another option could be leveraging the workload of software in production, e.g., through execution traces and logs, to understand which code components of the system are more heavily invoked and therefore claim for higher priority in performance testing.

## 6. Threats to validity

**Construct validity.** We focused solely on *Java software systems*. Our results may not generalize to systems developed in other programming languages. Nevertheless, Java is still among the most used programming languages<sup>8</sup>.

We restricted the coverage analysis to JMH microbenchmarks and JUnit tests since they are mature and widely adopted frameworks for developing performance and functional testing, respectively. Moreover, both these frameworks operate at the fine-grained level, as they are both used to test individual methods within a codebase. This commonality provides a fair basis for comparing their test coverage.

Using method-level coverage may have limitations, since this metric does not account for cases where performance/functional tests only partially cover

---

<sup>8</sup>Stack Overflow Developer Survey, <https://survey.stackoverflow.co/2023>.

the method statements. Our study results may change when employing a statement-level coverage metric. The decision to use method-level coverage stems from the significant technical challenges encountered in integrating JMH with traditional statement-level coverage tools, such as JaCoCo and Cobertura. To obtain statement-level coverage information, these tools modify the Java bytecode, which we observed could interfere with the execution of JMH microbenchmarks. Given these challenges, method-level coverage was deemed a reasonable compromise between the practicality of the study and the representativeness of the results. Furthermore, method-level coverage has been extensively employed in software performance research [58, 70, 53].

**External validity.** The presented analysis is limited to 28 open-source software systems. The findings may not be broadly generalizable; nonetheless, the selected systems are all well-known Java systems encompassing different domains (e.g., database systems, logging frameworks, and web servers). This limited number of subject systems is also motivated by an effort-intensive data collection, which required months of work (in multiple iterations) to get reliable results. This is a known issue in performance engineering that typically restricts the number of subject systems in empirical studies. Nevertheless, the number of subject systems used in our study is larger than most of the recent empirical studies on performance (e.g., see [30, 52, 11, 58, 12]).

**Internal validity.** We used a sampling-based CPU profiler to identify the methods covered by tests. These profilers operate by periodically capturing a program’s call stack during its execution. A limitation of this approach is the potential omission of call stack information. Since sampling is done at discrete intervals, short-lived function calls or those that fall between sampling points might not be captured. To mitigate this threat, we used `async-profiler`, which (to our knowledge) provides the lowest sampling rate for profiling Java software.

We employed various pre-processing techniques, including feature selection and class rebalancing, followed by hyperparameter tuning. However, other combinations of data preprocessing methods, algorithms, and hyperparameters could yield better results. Additionally, the source code features used for classifying methods might only partially capture the proper suitability of a method for performance testing. Nevertheless, we have used widely recognized source code features, which have been effectively utilized in previous software performance studies [60]. It is also important to note that, due

to the inherent variability in machine learning model training, the results may vary slightly if the experiments are repeated.

## 7. Related work

**Performance Testing.** The study most closely related to our work is that of Laaber and Leitner [30], which proposes a performance test quality metric inspired by mutation testing score, namely API benchmarking score (ABS). ABS is related to the concept of test coverage, as it represents the capability of the performance testing suite to find slowdowns. While Laaber and Leitner focus on defining a novel metric for test quality, our research evaluates the quality of existing performance testing suites using traditional code coverage metrics. Traini *et al.* [53] show that code components covered by performance tests tend to be less susceptible to refactoring. The execution time cost of performance testing is also related to this work. Researchers proposed approaches to reduce the time of performance testing without sacrificing results quality [52, 71, 72, 67].

**Test Coverage.** In [9], the authors describe Google’s code coverage infrastructure and how the computed information is visualized and used. The study demonstrates that most of the projects contain few unit tests, despite the opposite perception of the developers. The authors in [73] analyzed test coverage data on several widely used Python projects. The main finding is that the coverage strongly depends on the control flow structure. Moreover, the authors found that error-handling code is also neglected. In [8], the authors examine the question of coverage criteria as suite quality predictors from the perspective of non-researcher audience. Alves *et al.* [38] conceived an approach for estimating code coverage through static analysis, particularly slicing call graphs.

**Performance Bugs.** Jin *et al.* [74] empirically studied 110 real-world performance bugs collected from 5 open-source software repositories. A more extensive study was recently conducted by Zhao *et al.* [75], which investigated 570 performance issues from 13 open-source projects. Other empirical studies have focused on more specific domains, such as internet browsers [76], mobile applications [77], and JavaScript applications [64]. While our empirical findings may not directly correlate with the ability to uncover performance issues, there are strong indications of the relevance of code coverage for the efficacy of performance testing. Batch *et al.* [78] found that source

code covered by functional/performance tests is less prone to bugs. Ding *et al.* [11] showed that the code coverage of tests (i.e., scope) influences their capability to uncover performance bugs. The study of Jangali *et al.* [12] suggests that the extension of code coverage in performance testing improves the capability of detecting performance issues. These works indicate that incorporating code coverage analysis into performance testing can be beneficial for software performance assurance.

***Machine Learning for Performance Testing.*** Laaber et al. [60] studied the prediction of unstable software benchmarks using static source code features. They focused on benchmarks in open-source projects using the Go programming language. The study demonstrates the effectiveness of combining multiple static source code features and utilizing machine learning models to predict benchmark stability. Chen et al. [58] presented an approach to build prediction models to detect the tests that are prone to performance regression. Their work focuses on early identifying performance issues to mitigate the performance regressions introduced by new code commits at the test level. They exploit the historical commit information to extract features. Machine learning models are employed to predict performance regressions immediately after a new commit is introduced. The study of Zhao et al. [79] proposes an approach for enhancing performance bug prediction using performance code metrics. The study highlights the effectiveness of combining various performance code metrics with machine learning techniques to identify performance issues. They integrate code characteristics of performance bugs with traditional source code metrics to build machine learning models for predicting presence of performance bugs in the source code.

## 8. Conclusion and future work

This paper presented a comprehensive empirical study focused on performance testing coverage. Our findings revealed the limited coverage of current performance testing suites and the significant execution time cost associated with them. The results of this work suggest opportunities to enhance the coverage of performance testing suites, by emphasizing the necessity to enlighten practitioners about these prevalent limitations.

We have intentionally considered in this paper the concept of code coverage that usually relates to functional testing. Additional metrics should be considered for a sharper concept of performance test coverage, like workload

and operational profile. However, these metrics are quite difficult to collect and may sensibly vary for the same application in different contexts. Therefore, we have intended to explore the extent at which performance testing can be solely based on code coverage. Additionally, we have investigated the relationship between source code features and the likelihood of performance testing opportunity. Our goal was to identify recurring characteristics in the source code that could drive performance testing efforts.

As suggested by our findings, the real-world adoption of performance testing techniques might be hampered by their substantial execution time costs. The evidence provided in this paper sustains the idea that the limited coverage of performance tests (as compared to functional ones) only stems from technical issues (e.g., time limits, problems to collect dynamic metrics). Indeed, a wider coverage is desirable as it would allow to identify performance issues in code sections that, for the above reasons, are usually not considered.

To address this challenge, researchers can leverage the automated generation of performance tests [12, 18]. One promising direction is the development of “smart” test selection strategies that can reduce the execution time of a performance testing suite without compromising its effectiveness. The selection of code components for performance testing appears to be driven primarily by objectives that cannot be inferred solely from source code features. Instead, these objectives may be closely tied to domain-specific reasons or associated with the operational profile of the software system. Therefore, future performance test generation techniques should better consider these aspects rather than focusing solely on source code features. We encourage future research efforts directed to address this challenge.

## References

- [1] J. Brutlag, Google ai blog: Speed matters (2009).  
URL <https://ai.googleblog.com/2009/06/speed-matters.html>
- [2] S. Olenski, Why brands are fighting over milliseconds (2016).  
URL <https://www.forbes.com/sites/steveolenski/2016/11/10/why-brands-are-fighting-over-milliseconds/>
- [3] E. J. Weyuker, Experience with performance testing of software systems: issues, an approach, and case study, *IEEE Transactions on Software Engineering* 26 (12) (2000) 1147 – 1156. doi:10.1109/32.888628.

- [4] W. Behutiye, P. Karhapää, L. López, X. Burgués, S. Martínez-Fernández, A. M. Vollmer, P. Rodríguez, X. Franch, M. Oivo, Management of quality requirements in agile and rapid software development: A systematic mapping study, *Information and Software Technology* 123 (2020) 106225.
- [5] L. Traini, Exploring Performance Assurance Practices and Challenges in Agile Software Development: An Ethnographic Study, *Empirical Software Engineering* 27 (3) (2022) 74. doi:10.1007/s10664-021-10069-3. URL <https://doi.org/10.1007/s10664-021-10069-3>
- [6] W. Alsaqaf, M. Daneva, R. Wieringa, Quality requirements challenges in the context of large-scale distributed agile: An empirical study, *Information and Software Technology* 110 (2019) 39 – 55.
- [7] L. Inozemtseva, R. Holmes, Coverage is not strongly correlated with test suite effectiveness, in: *Proceedings of the 36th International Conference on Software Engineering, ICSE 2014, ACM, 2014*, p. 435–445.
- [8] R. Gopinath, C. Jensen, A. Groce, Code coverage for suite evaluation by developers, in: *Proceedings of the 36th International Conference on Software Engineering, ICSE 2014, ACM, 2014*, p. 72–82.
- [9] M. Ivanković, G. Petrović, R. Just, G. Fraser, Code coverage at google, in: *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2019, ACM, 2019*, p. 955–963.
- [10] P. Piwowarski, M. Ohba, J. Caruso, Coverage measurement experience during function test, in: *Proceedings of the 15th International Conference on Software Engineering, ICSE '93, IEEE Computer Society Press, 1993*, p. 287–301.
- [11] Z. Ding, J. Chen, W. Shang, Towards the use of the readily available tests from the release pipeline as performance tests: are we there yet?, in: G. Rothermel, D. Bae (Eds.), *ICSE '20: 42nd International Conference on Software Engineering, Seoul, South Korea, 27 June - 19 July, 2020, ACM, 2020*, pp. 1435–1446. doi:10.1145/3377811.3380351. URL <https://doi.org/10.1145/3377811.3380351>

- [12] M. Jangali, Y. Tang, N. Alexandersson, P. Leitner, J. Yang, W. Shang, Automated generation and evaluation of jmh microbenchmark suites from unit tests, *IEEE Transactions on Software Engineering* 49 (4) (2023) 1704–1725. doi:10.1109/TSE.2022.3188005.
- [13] M. Imran, V. Cortellessa, D. Di Ruscio, R. Rubei, L. Traini, An empirical study on code coverage of performance testing, in: *Proceedings of the 28th International Conference on Evaluation and Assessment in Software Engineering, EASE '24*, Association for Computing Machinery, New York, NY, USA, 2024, p. 48–57. doi:10.1145/3661167.3661196. URL <https://doi.org/10.1145/3661167.3661196>
- [14] S. Yoo, M. Harman, Regression testing minimization, selection and prioritization: A survey, *Softw. Test. Verif. Reliab.* 22 (2) (2012) 67–120.
- [15] C. Laaber, T. Yue, S. Ali, Evaluating search-based software microbenchmark prioritization, *IEEE Transactions on Software Engineering* (2024) 1–16doi:10.1109/TSE.2024.3380836.
- [16] C. Laaber, H. C. Gall, P. Leitner, Applying test case prioritization to software microbenchmarks, *Empirical Software Engineering* 26 (6) (2021) 133. doi:10.1007/s10664-021-10037-x. URL <https://doi.org/10.1007/s10664-021-10037-x>
- [17] A. B. De Oliveira, S. Fischmeister, A. Diwan, M. Hauswirth, P. F. Sweeney, Perphecy: Performance regression test selection made simple but effective, in: *2017 IEEE International Conference on Software Testing, Verification and Validation (ICST)*, 2017, pp. 103–113. doi:10.1109/ICST.2017.17.
- [18] M. Rodriguez-Cancio, B. Combemale, B. Baudry, Automatic microbenchmark generation to prevent dead code elimination and constant folding, in: *Proceedings of the 31st IEEE/ACM International Conference on Automated Software Engineering, ASE '16*, ACM, 2016, p. 132–143. doi:10.1145/2970276.2970346. URL <https://doi.org/10.1145/2970276.2970346>
- [19] M. Imran, V. Cortellessa, D. Di Ruscio, R. Rubei, L. Traini, Is Code Coverage of Performance Tests Related to Source Code Features? An Empirical Study on Open-Source Systems - Replication Package (2024).

URL <https://github.com/SpencerLabAQ/performance-test-coverage>

- [20] C. U. Smith, L. G. Williams, Performance solutions: a practical guide to creating responsive, scalable software, Vol. 23, Addison-Wesley Reading, 2002.
- [21] L. Chung, B. A. Nixon, E. Yu, J. Mylopoulos, Non-functional requirements in software engineering, Vol. 5, Springer Science & Business Media, 2012.
- [22] ISO/IEC, ISO/IEC 25010 - System and software quality models, <https://iso25000.com/index.php/en/iso-25000-standards/iso-25010>, retrieved July, 2024.
- [23] Y. Zhao, L. Xiao, W. Xiao, B. Chen, Y. Liu, Localized or architectural: an empirical study of performance issues dichotomy, in: 2019 IEEE/ACM 41st International Conference on Software Engineering: Companion Proceedings (ICSE-Companion), IEEE, 2019, pp. 316–317.
- [24] F. I. Vokolos, E. J. Weyuker, Performance testing of software systems, in: Proceedings of the 1st International Workshop on Software and Performance, 1998, pp. 80–87.
- [25] E. J. Weyuker, F. I. Vokolos, Experience with performance testing of software systems: issues, an approach, and case study, IEEE transactions on software engineering 26 (12) (2000) 1147–1156.
- [26] S. He, T. Liu, P. Lama, J. Lee, I. K. Kim, W. Wang, Performance testing for cloud computing with dependent data bootstrapping, in: 2021 36th IEEE/ACM International Conference on Automated Software Engineering (ASE), IEEE, 2021, pp. 666–678.
- [27] S. He, G. Manns, J. Saunders, W. Wang, L. Pollock, M. L. Soffa, A statistics-based performance testing methodology for cloud applications, in: Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, 2019, pp. 188–199.

- [28] Z. M. Jiang, A. E. Hassan, A survey on load testing of large-scale software systems, *IEEE Transactions on Software Engineering* 41 (11) (2015) 1091–1118.
- [29] D. Eadline, Micro-benchmarks vs macro-benchmarks, accessed: July, 2024 (September 2005).  
URL <https://www.clustermonkey.net/Benchmarking-Methods/micro-benchmarks-vs-macro-benchmarks.html>
- [30] C. Laaber, P. Leitner, An evaluation of open-source software microbenchmark suites for continuous performance assessment, *ACM*, 2018.
- [31] A. Shipilev, Nanotrusting nanotime, <https://shipilev.net/blog/2014/nanotrusting-nanotime/>, accessed: 2024-07-04 (2014).
- [32] D. Costa, C.-P. Bezemer, P. Leitner, A. Andrzejak, What’s wrong with my benchmark results? studying bad practices in jmh benchmarks, *IEEE Transactions on Software Engineering* 47 (7) (2019) 1452–1467.
- [33] T. Kalibera, R. Jones, Quantifying performance changes with effect size confidence intervals, *arXiv preprint arXiv:2007.10899* (2020).
- [34] P. Stefan, V. Horky, L. Bulej, P. Tuma, Unit testing performance in java projects: Are we there yet?, in: *Proceedings of the 8th ACM/SPEC on International Conference on Performance Engineering*, 2017, pp. 401–412.
- [35] Google, Google Caliper (Year of the last commit or latest release).  
URL <https://github.com/google/caliper>
- [36] N. O’Connor, JUnitPerf (Year of the last commit or latest release).  
URL <https://github.com/noconnor/JUnitPerf>
- [37] O. Corporation, Java Microbenchmarking Harness (JMH) (2016).  
URL <https://openjdk.org/projects/code-tools/jmh/>
- [38] T. L. Alves, J. Visser, Static estimation of test coverage, in: *2009 Ninth IEEE International Working Conference on Source Code Analysis and Manipulation*, IEEE, 2009, pp. 55–64.

- [39] A. Rohella, S. Takada, Testing android applications using multi-objective evolutionary algorithms with a stopping criteria., in: SEKE, 2018, pp. 308–307.
- [40] Z. Zhou, Y. Zhou, C. Fang, Z. Chen, Y. Tang, Selectively combining multiple coverage goals in search-based unit test generation, in: Proceedings of the 37th IEEE/ACM International Conference on Automated Software Engineering, 2022, pp. 1–12.
- [41] C. Häubl, C. Wimmer, H. Mössenböck, Deriving code coverage information from profiling data recorded for a trace-based just-in-time compiler, in: Proceedings of the 2013 International Conference on Principles and Practices of Programming on the Java Platform: Virtual Machines, Languages, and Tools, 2013, pp. 1–12.
- [42] Q. Yang, J. J. Li, D. Weiss, A survey of coverage based testing tools, in: Proceedings of the 2006 international workshop on Automation of software test, 2006, pp. 99–103.
- [43] M. M. U. Alam, T. Liu, G. Zeng, A. Muzahid, Syncperf: Categorizing, detecting, and diagnosing synchronization performance bugs, in: Proceedings of the Twelfth European Conference on Computer Systems, 2017, pp. 298–313.
- [44] C. Laaber, H. C. Gall, P. Leitner, Applying test case prioritization to software microbenchmarks, *Empirical Software Engineering* 26 (6) (2021) 133.
- [45] S. Brown, A. Mitchell, J. F. Power, A coverage analysis of java benchmark suites (2005).
- [46] P. Leitner, C.-P. Bezemer, An exploratory study of the state of practice of performance testing in java-based open source projects, in: Proceedings of the 8th ACM/SPEC on International Conference on Performance Engineering, ICPE '17, ACM, 2017, p. 373–384. doi: 10.1145/3030207.3030213.  
URL <https://doi.org/10.1145/3030207.3030213>
- [47] N. V. Chawla, K. W. Bowyer, L. O. Hall, W. P. Kegelmeyer, Smote: synthetic minority over-sampling technique, *Journal of artificial intelligence research* 16 (2002) 321–357.

- [48] G. E. Batista, R. C. Prati, M. C. Monard, A study of the behavior of several methods for balancing machine learning training data, *ACM SIGKDD explorations newsletter* 6 (1) (2004) 20–29.
- [49] I. Guyon, J. Weston, S. Barnhill, V. Vapnik, Gene selection for cancer classification using support vector machines, *Machine learning* 46 (2002) 389–422.
- [50] J. Jiarpakdee, C. Tantithamthavorn, C. Treude, Autospearman: Automatically mitigating correlated software metrics for interpreting defect models, in: *2018 IEEE International Conference on Software Maintenance and Evolution (ICSME)*, IEEE Computer Society, 2018, pp. 92–103.
- [51] M. Avinash, M. Nithya, S. Aravind, Automated machine learning-algorithm selection with fine-tuned parameters, in: *2022 6th International Conference on Intelligent Computing and Control Systems (ICICCS)*, IEEE, 2022, pp. 1175–1180.
- [52] C. Laaber, S. Würsten, H. C. Gall, P. Leitner, Dynamically reconfiguring software microbenchmarks: Reducing execution time without sacrificing result quality, in: *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2020*, ACM, 2020, p. 989–1001.
- [53] L. Traini, D. Di Pompeo, M. Tucci, B. Lin, S. Scalabrino, G. Bavota, M. Lanza, R. Oliveto, V. Cortellessa, How software refactoring impacts execution time, *ACM Trans. Softw. Eng. Methodol.* 31 (2) (dec 2021). doi:10.1145/3485136.  
URL <https://doi.org/10.1145/3485136>
- [54] L. Traini, V. Cortellessa, D. Di Pompeo, M. Tucci, Towards effective assessment of steady state performance in java software: are we there yet?, *Empirical Software Engineering* 28 (1) (2022) 13. doi:10.1007/s10664-022-10247-x.  
URL <https://doi.org/10.1007/s10664-022-10247-x>
- [55] M. L. Collard, M. J. Decker, J. I. Maletic, Lightweight transformation and fact extraction with the srml toolkit, in: *2011 IEEE 11th Interna-*

- tional Working Conference on Source Code Analysis and Manipulation, 2011, pp. 173–184. doi:10.1109/SCAM.2011.19.
- [56] A. Walker, M. Coffey, P. Tisnovsky, T. Cerny, On limitations of modern static analysis tools, in: K. J. Kim, H.-Y. Kim (Eds.), *Information Science and Applications*, Springer Singapore, 2020, pp. 577–586.
  - [57] T. J. McCabe, A complexity measure, *IEEE Transactions on software Engineering* (4) (1976) 308–320.
  - [58] J. Chen, W. Shang, E. Shihab, Perfjit: Test-level just-in-time prediction for performance regression introducing commits, *IEEE Transactions on Software Engineering* 48 (5) (2022) 1529–1544.
  - [59] J. P. S. Alcocer, A. Bergel, M. T. Valente, Prioritizing versions for performance regression testing: the pharo case, *Science of Computer Programming* 191 (2020) 102415.
  - [60] C. Laaber, M. Basmaci, P. Salza, Predicting unstable software benchmarks using static source code features, *Empirical Software Engineering* 26 (6) (2021) 114.
  - [61] P. Huang, X. Ma, D. Shen, Y. Zhou, Performance regression testing target prioritization via performance risk analysis, in: *Proceedings of the 36th International Conference on Software Engineering*, 2014, pp. 60–71.
  - [62] G. Jin, L. Song, X. Shi, J. Scherpelz, S. Lu, Understanding and detecting real-world performance bugs, *ACM SIGPLAN Notices* 47 (6) (2012) 77–88.
  - [63] J. P. S. Alcocer, A. Bergel, Tracking down performance variation against source code evolution, *ACM SIGPLAN Notices* 51 (2) (2015) 129–139.
  - [64] M. Selakovic, M. Pradel, Performance issues and optimizations in javascript: an empirical study, in: *Proceedings of the 38th International Conference on Software Engineering*, 2016, pp. 61–72.
  - [65] Y. Zhao, L. Xiao, X. Wang, L. Sun, B. Chen, Y. Liu, A. B. Bondi, How are performance issues caused and resolved?-an empirical study from a design perspective, in: *Proceedings of the ACM/SPEC International Conference on Performance Engineering*, 2020, pp. 181–192.

- [66] D. Costa, A. Andrzejak, J. Seboek, D. Lo, Empirical study of usage and performance of java collections, in: Proceedings of the 8th ACM/SPEC on International Conference on Performance Engineering, 2017, pp. 389–400.
- [67] T. Kalibera, R. Jones, Rigorous benchmarking in reasonable time, in: Proceedings of the 2013 International Symposium on Memory Management, ISMM '13, ACM, 2013, p. 63–74. doi:10.1145/2491894.2464160.  
URL <https://doi.org/10.1145/2491894.2464160>
- [68] K. H. Brodersen, C. S. Ong, K. E. Stephan, J. M. Buhmann, The balanced accuracy and its posterior distribution, in: 2010 20th International Conference on Pattern Recognition, 2010, pp. 3121–3124. doi:10.1109/ICPR.2010.764.
- [69] H. Samoaa, P. Leitner, An exploratory study of the impact of parameterization on jmh measurement results in open-source projects, in: Proceedings of the ACM/SPEC International Conference on Performance Engineering, ICPE '21, ACM, 2021, p. 213–224. doi:10.1145/3427921.3450243.  
URL <https://doi.org/10.1145/3427921.3450243>
- [70] J. Chen, Z. Ding, Y. Tang, M. Sayagh, H. Li, B. Adams, W. Shang, Iopv: On inconsistent option performance variations, in: Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2023, ACM, 2023, p. 845–857. doi:10.1145/3611643.3616319.  
URL <https://doi.org/10.1145/3611643.3616319>
- [71] H. M. Alghmadi, M. D. Syer, W. Shang, A. E. Hassan, An automated approach for recommending when to stop performance tests, in: 2016 IEEE International Conference on Software Maintenance and Evolution (ICSME), 2016, pp. 279–289. doi:10.1109/ICSME.2016.46.
- [72] S. He, G. Manns, J. Saunders, W. Wang, L. Pollock, M. L. Soffa, A statistics-based performance testing methodology for cloud applications, in: Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2019, ACM, 2019, p. 188–199.

- [73] H. Zhai, C. Casalnuovo, P. Devanbu, Test coverage in python programs, in: 2019 IEEE/ACM 16th International Conference on Mining Software Repositories (MSR), IEEE, 2019, pp. 116–120.
- [74] G. Jin, L. Song, X. Shi, J. Scherpelz, S. Lu, Understanding and detecting real-world performance bugs, in: Proceedings of the 33rd ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '12, ACM, 2012, p. 77–88. doi:10.1145/2254064.2254075. URL <https://doi.org/10.1145/2254064.2254075>
- [75] Y. Zhao, L. Xiao, A. B. Bondi, B. Chen, Y. Liu, A large-scale empirical study of real-life performance issues in open source projects, IEEE Transactions on Software Engineering 49 (2) (2023) 924–946.
- [76] S. Zaman, B. Adams, A. E. Hassan, A qualitative study on performance bugs, in: Proceedings of the 9th IEEE Working Conference on Mining Software Repositories, MSR '12, IEEE Press, 2012, p. 199–208.
- [77] A. Mazuera-Rozo, C. Trubiani, M. Linares-Vásquez, G. Bavota, Investigating types and survivability of performance bugs in mobile apps, Empirical Software Engineering 25 (2020) 1644–1686.
- [78] T. Bach, A. Andrzejak, R. Pannemans, D. Lo, The impact of coverage on bug density in a large industrial software project, in: 2017 ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM), IEEE, 2017, pp. 307–313.
- [79] G. Zhao, S. Georgiou, S. Hassan, Y. Zou, D. Truong, T. Corbin, Enhancing performance bug prediction using performance code metrics, in: Proceedings of the 21st International Conference on Mining Software Repositories, MSR '24, Association for Computing Machinery, New York, NY, USA, 2024, p. 50–62. doi:10.1145/3643991.3644920. URL <https://doi.org/10.1145/3643991.3644920>