# Predicting post-release defects with knowledge units (KUs) of programming languages: an empirical study

**Md Ahasanuzzaman · Gustavo A. Oliva · Ahmed E. Hassan · Zhen Ming (Jack) Jiang**

**Abstract** Defect prediction plays a crucial role in software engineering, enabling developers to identify defect-prone code and improve software quality. While extensive research has focused on refining machine learning models for defect prediction, the exploration of new data sources for feature engineering remains limited. Defect prediction models primarily rely on traditional metrics such as product, process, and code ownership metrics, which, while effective, do not capture language-specific traits that may influence defect proneness. To address this gap, we introduce Knowledge Units (KUs) of programming languages as a novel feature set for analyzing software systems and defect prediction. A KU is a cohesive set of key capabilities that are offered by one or more building blocks of a given programming language. We conduct an empirical study leveraging 28 KUs that are derived from Java certification exams and compare their effectiveness against traditional metrics in predicting post-release defects across 28 releases of 8 well-maintained Java software systems. Our results show that KUs provide significant predictive power, achieving a median AUC of 0.82, outperforming individual group of traditional metric-based models (e.g., process, product and ownership metrics). Among KU features, Method & Encapsulation, Inheritance, and Exception Handling emerge as the most influential predictors. Furthermore, combining KUs with traditional metrics enhances prediction performance, yielding a median AUC of 0.89. We also introduce a cost-effective model using only 10 features (5 KUs and 5 traditional metrics), which maintains strong predictive performance while reducing feature engineering costs. Our findings demonstrate the value of KUs in predicting post-release defects, offering a complementary perspec-

✉ Md Ahasanuzzaman, Gustavo A. Oliva, and Ahmed E. Hassan
Software Analysis and Intelligence Lab (SAIL), School of Computing
Queen's University, Kingston, Ontario, Canada
E-mail: 16ma87@queensu.ca,{gustavo,ahmed}@cs.queensu.ca
✉ Zhen Ming (Jack) Jiang
York University, Toronto, Ontario, Canada
E-mail: zmjiang@cse.yorku.ca

tive to traditional metrics. This study can be helpful to researchers who wish to analyze software systems from a perspective that is complementary to that of traditional metrics.

**Keywords** Defect prediction , Software metrics, code metrics, Java, certification exam, software engineering, and machine learning

## 1 Introduction

Defect prediction has been a long-standing area of research in the field of software engineering, aimed at identifying potential defect-prone code in software systems [27, 34, 38, 40, 48]. While the defect prediction field has seen substantial research in developing new machine learning approaches (e.g., bootstrapped by the recent advances and popularization of deep learning techniques [7, 21, 31, 57], the exploration of novel types of data for feature engineering remains less explored. The continual improvement in the defect prediction field demands innovation not just in modeling techniques but also in the data used for feature engineering. The more indicators researchers have, the better equipped they are to understand defects and hypothesize their underlying causes.

Product, process and code ownership metrics have been extensively used for defect prediction over several decades. We collectively refer to them as traditional metrics. These metrics characterize the size (e.g., LOC), complexity (e.g., cyclomatic complexity), and structure (e.g., the CK suite [8]) of object-oriented software systems. Most importantly, these metrics are *programming-language agnostic*, in the sense that their definitions are applicable to and computable for any object-oriented system (irrespective of the used programming language). As an inherent limitation, traditional metrics do not reveal system *traits* that are tied to certain building blocks of a given programming language. As an illustrative example, let us consider a system written in Java and that makes use of the Java Concurrency API (i.e., a few classes from the `java.util.concurrent` package). Detecting the use of such an API reveals that concurrent code is a *trait* of that system. Such a trait is relevant: specific defects might arise because of it, not every developer is proficient in concurrent programming, and debugging a concurrent program is often more difficult than debugging a single-threaded one. Yet, traditional metrics would not capture this concurrency trait.

Towards filling this gap, this paper reports an empirical study on the usage of knowledge units (KUs) of the Java programming language for predicting post-release defects. As defined in our prior study [2], a KU is a *cohesive set of key capabilities that are offered by one or more building blocks of a given programming language* (see Figure 1). The word "capabilities" refer to the main things that a programmer can do with a certain building block of a programming language. For instance, using the *Concurrency API* building block, a developer can create worker threads to execute tasks concurrently. Therefore, it is reasonable to assume that Java has a *Concurrency KU*, which includes

a cohesive set of key concurrent program capabilities offered by the Concurrency API (building block). Our hypothesis is that by capturing system traits such as *concurrency* through KUs, we will be able to design more accurate post-release defect prediction models compared to those relying exclusively on traditional metrics.

We note that KUs are also able to capture more nuanced code characteristics (e.g., compared to concurrent programming) that may still influence defect-proneness. For instance, a *method* is a fundamental building block in Java. Developers can create methods with parameters, overloaded methods, methods with variable-length arguments, and define access modifiers for methods. Not every object-oriented language supports these capabilities. Therefore, it is also reasonable to assume that Java has a *Method & Encapsulation KU*, which captures the key capabilities associated with methods and constructors. KUs capture deeper concepts than merely counting keywords. For instance, detecting the Method & Encapsulation KU involves tracking the usage of variable length arguments in methods and the application of overloaded constructors and constructor chaining. To achieve this, we use the Eclipse JDT framework to parse the code's Abstract Syntax Tree (AST), allowing us to gather nuanced and richer information that goes much beyond simple keyword matching. By incorporating such granular details, KUs offer a more comprehensive perspective on software structure that extends beyond traditional metrics.

As described in our prior work [2], we elicit KUs by manually analyzing three certification exams for the Java programming language (Section 2.2): Java SE 8 Programmer I Exam, Java SE 8 Programmer II Exam, and Java EE Developer Exam. We elicit a total of 28 KUs, which cover a wide range of capabilities (A). To count the occurrences of each KU in a Java source code file, we built a custom KU detector on top of the Eclipse JDT framework (Section 2.3). We refer to the set of metrics corresponding to the counts of each KU as *KU metrics*.

Our data collection is designed as follows. First, we select the defect dataset that was curated by Yatish et al. [56]. This dataset contains post-release defect data for 28 releases that originate from 8 Java systems. For each source code file in a release, we compute our 28 KU metrics. Additionally, we reuse all 65 traditional metrics provided in the defect dataset, including 54 product metrics, 5 process metrics, and 6 code ownership metrics for each source code file. Finally, we build defect prediction models using these two sets of metrics to address the following research questions (RQs):

**RQ1: How effective are KU features for predicting post-release defects?** First, we use the KU metrics to build a defect prediction model for each of the studied releases. We refer to this model as KUM. To judge the performance of KUM, we build four baseline models using traditional metrics: (i) PROD (which is built with product metrics), (ii) PROC (which is built with process metrics), (iii) OWN (which is built with code ownership metrics) and (iv) TM (which is built using all product, process and code ownership metrics). To compare the performance of KUM and any of the baseline

models, we perform a statistical test using a two-sided Wilcoxon signed-rank test. Additionally, we employ Cliff's delta [33] as an effect size measure to assess the practical significance of the differences between the distributions. We also rank the studied models using the Scott-Knott ESD (SK-ESD) technique [18, 36, 55]. Our results indicate that:

*KUM achieves a median AUC of 0.82 (a defect prediction model with an AUC of 0.70 or higher is typically considered suitable for practical use). KUM outperforms three baseline models: PROD, PROC and OWN. The normalized AUC improvement of KUM over PROD ranges from 3.6% to 22.2% (median 11.4%), improvement over PROC ranges from 6.7% to 68.9% (median 42.4%), and improvement over OWN ranges from 3.3% to 71.9% (median (33.3%). That is, KUM not only outperforms a model built on the same raw data (PROD), but also those built on process (PROC) and ownership (OWN) data. TM is the top-performing model, achieving a median AUC of 0.85*

**RQ2: What are the most important KU features for predicting post-release defects?** To determine the most important features of KUM, we utilize SHAP (SHapley Additive exPlanation). SHAP is a robust, flexible, and widely used method for model interpretation [37, 49]. Next, we apply the SK-ESD technique on the distribution of SHAP values to rank features. Our results indicate that:

*Method & Encapsulation, Inheritance, and Exception are identified as the top three most important KU features of KUM for predicting post-release defects.*

**RQ3: Can KUM be made more accurate by combining its features with those from TM?** Based on the encouraging results from RQ1, we join all the features of KUM and TM to build a combined model, which we call KUM+TM. We compare the performance of KUM+TM with that of KUM and TM using the same approach from RQ1. To investigate the importance of feature dimensions in KUM+TM, we employ the same SHAP approach from RQ2. Our results indicate that:

*KUM+TM achieves a median AUC of 0.89 and outperforms both KUM and TM across all of the studied releases. The normalized AUC improvement of KUM+TM over TM ranges from 11.1% to 44.4% (median 22.5%) and the improvement over KUM ranges from 10.5% to 72.2% (median 31.5%) across the studied releases. The most important feature is the number of active developers (ADEV, a process feature). The top ranked KU feature is Method & Encapsulation at rank 4. Inheritance and Exception Handling KU features are also among the top ten most important features.*

**RQ4: Can a cost-effective model be created by combining specific features from KUM and TM?** A cost-effective model that uses fewer features and yet maintains decent performance is valuable, as it reduces operational costs associated with extensive feature engineering. Towards this goal, we build a cost-effective combined model, which we call COST_EFF. This model is built using the top five features of KUM and the top five traditional metric features of TM. This results in a combined model with only 10 features,

compared to the 65 features in TM and the 28 features in KULTC. Our results indicate that:

*The COST_EFF achieves a median AUC of 0.87 and outperforms TM. The normalized AUC improvement of COST_EFF over TM ranges from 3.7% to 31.3% (median 11.1%). Therefore, we improve over the model built with all traditional metrics while still keeping feature engineering cost at a much more reasonable level.*

The main contribution of this paper is to evaluate the effectiveness of KUs of the Java programming language for predicting post-release defects. We believe that our findings can be helpful and encouraging to researchers who wish to analyze software systems from a perspective that is complementary to traditional metrics. Our encouraging results indicate that further refining the concept of KUs and exploring alternative elicitation techniques can potentially lead to even higher-performing defect prediction models. To bootstrap future work in this area, we provide online supplementary material with the data analyzed as part of this study[1].

**Paper organization.** Section 2 defines KUs and presents our approach for detecting them. Section 3 describes our data collection approach. Section 4 presents the motivation, approach, and findings that are associated with our research questions. Section 5 discusses the trade-offs between traditional metrics and KUs for defect prediction and future directions to expanding KUs. Section 6 describes the threats to the validity of our findings. Section 7 discusses related work. Finally, Section 8 concludes the paper.

## 2 Knowledge Units (KUs)

We introduce the concept of knowledge units (KUs) of programming language in our prior work [2]. To make this paper as self-contained as possible, in the following we provide a brief introduction to KUs (Section 2.1) and discuss our approaches for eliciting (Section 2.2) and detecting them (Section 2.3).

### 2.1 Definition

Every programming language has its own building blocks. The building blocks of a programming language are the "blocks" that developers use to "build" (write) code. In the case of Java, we consider them to be fundamental language constructs (e.g., arrays, if/else statements, try/catch statements) and APIs (e.g., Generic and Collection API, Concurrency API and String API). Each building block offers a *set of capabilities* essentially representing the actions a developer can perform using that building block. We define a Knowledge Unit (KU) as a cohesive set of **key** capabilities that are offered by **one or more building blocks** of a given programming language. The inclusion of "key"
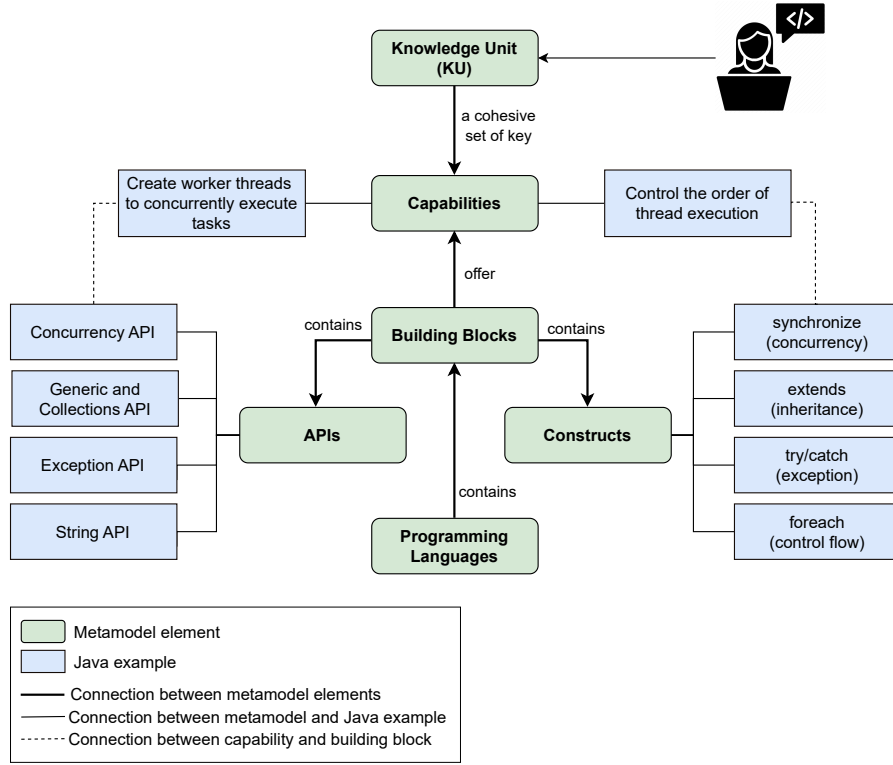
---

[1] http://bit.ly/3GoSmHL

**Fig. 1** Our metamodel for knowledge units (KUs).

in our definition aims to ensure that KUs are centered around fundamental capabilities rather than those that are overly specific. Figure 1 illustrates key capabilities that are associated with the `Concurrency API` (left-hand side) and the `synchronize` language construct (right-hand side).

2.2 Eliciting KUs

Certification exams of a programming language (e.g., Oracle Java SE and Java EE certification exams for Java) aim to determine the proficiency of a developer in using the key capabilities offered by the building blocks of that language. Therefore, we can say that *certification exams capture the KUs of a programming language.* For instance, one of the topics covered in the Java SE 8 Programmer II certification exam is "Generics and Collections." The subtopics of "Generics and Collections" that are covered in the exam include: (i) create and use a generic class, (ii) create and use `ArrayList`, `TreeSet`, `TreeMap`, and `ArrayDeque` objects, (iii) use `java.util.Comparator` and `java.lang.Comparable` interfaces, and (iv) iterate using `forEach` meth-

ods of a `List`. We interpret such a list of subtopics as the key capabilities that are offered by the *Generics and Collections* building block of the Java programming language. From this interpretation, we infer that Java has a *Generics and Collections* KU.

We reuse the KUs that we elicited in our prior work [2]. To elicit those KUs, we relied on the Oracle certification exams for the Java programming language. Oracle offers certification exams for different Java editions, such as Java Standard Edition (Java SE) and Java Enterprise Edition (Java EE). Specifically, for Java SE, there are two levels of certification exams:(i) *Java SE 8 Programmer I Certification exam* [46], and (ii) *Java SE 8 Programmer II Certification exam*[44]. For Java EE, only one version of the exam is offered by Oracle: *Oracle Certified Professional, Java EE Application Developer Certification exam* [45]. We elicited KUs directly from the topics outlined in these exams. In most cases, an exam topic was interpreted as a KU (e.g., the Generics and Collections topic of Java SE 8 Programmer II Certification exam was interpreted as "Generics and Collections") and its subtopics (e.g., create a generic class, use `TreeSet, TreeMap and ArrayList` and use `java.util.Comparator`) were interpreted as the key capabilities of that KU. At the end of this process, we elicited 28 KUs (Table 1). We refer to our prior work [2] for the detailed description of this mapping process.

## 2.3 Detection of KUs

Analogously to our prior work [2], we employ static analysis to detect KUs from a given Java project's release. First, we parse all Java source files of the release using the Eclipse JDT framework [11]. The JDT framework creates an Abstract Syntax Tree (AST) and provides visitors for every element of such tree, including *names*, *types* (class, interfaces), *expressions*, *statements*, and *declarations*.[2]
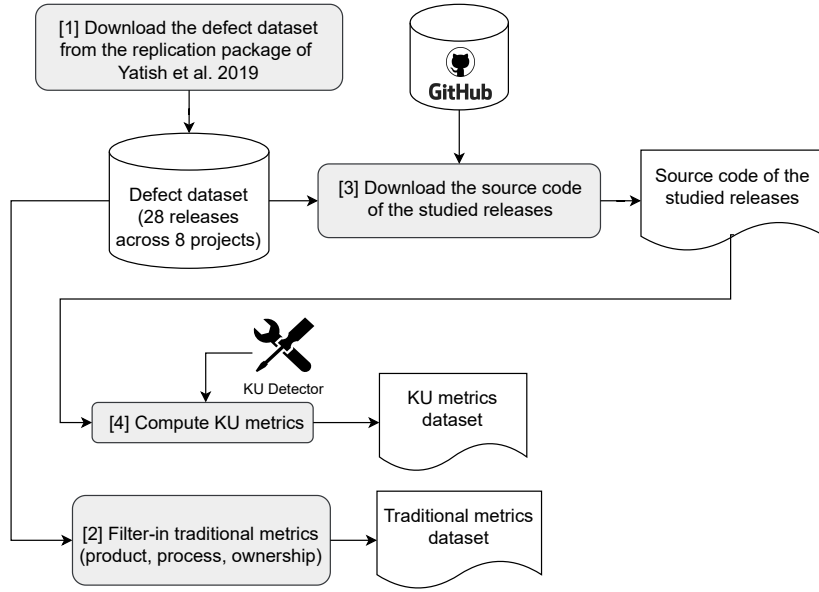
To detect a KU, we resort to the set of key capabilities of that unit (Tables 6). For instance, to detect the presence of the *Generics and Collections KU*, we use the visitors provided by JDT to detect the capabilities associated with such unit, namely: (i) create and use a generic class, (ii) create and use `ArrayList, TreeSet, TreeMap`, and `ArrayDeque` objects, and (iii) use `java.util.Comparator` and `java.lang.Comparable` interfaces. For instance, using type visitors, we can determine whether a given class is a generic class. Similarly, using statement visitors, we can find instantiations of an ArrayList. We also use JDT to collect binding information of methods, classes or variables (e.g., `org.eclipse.jdt.core.dom.IVariableBinding` resolves the binding information for a variable). In our study, we exclude type bindings that are related to third-party libraries since our primary interest lies in studying the KUs that are associated with the source code written by the developers of the studied projects.

---

[2] The grammar employed by JDT can be seen at `https://github.com/eclipse-jdt/eclipse.jdt.core/blob/master/org.eclipse.jdt.core.compiler.batch/grammar/java.g`

**Table 1** The list of knowledge units (KUs) elicited for the Java programming language.

| Knowledge unit (KU) | Definition |
| --- | --- |
| **[K1]** Data Type | The declaration and initialization of different types of variables (e.g., primitive type and parameterized type.) |
| **[K2]** Operator & Decision | The usage of different Java operators (e.g., assignment, logical, and bit-wise operators) and conditional statements (e.g, if, if-else, and switch statements). |
| **[K3]** Array | The declaration, instantiation, initialization and the usage of one-dimensional and multi-dimensional arrays. |
| **[K4]** Loop | The execution of a set of instruction-s/methods repeatedly using for, while, and do-while statements and the skipping and stopping of a repetitive execution of instructions and methods using continue and break statements. |
| **[K5]** Method & Encapsulation | The creation of methods with parameters, the use of overloaded methods and constructors, the usage of constructor chaining, the creation of methods with variable length arguments, and the usage of different access modifiers. This KU also describes encapsulation mechanisms, such as creating a set and get method for controlling data access, generating immutable classes, and updating object type parameters of a method. |
| **[K6]** Inheritance | Developing code with child class and parent class relationship, using polymorphism (e.g., developing code with overridden methods), creating abstract classes and interfaces, and accessing methods and fields of the parent's class. |
| **[K7]** Advanced Class Design | Developing code that uses the final keyword, creating inner classes including static inner classes, local classes, nested classes, and anonymous inner classes, using enumerated types including methods and constructors in an enum type, and developing code using @override annotator. |
| **[K8]** Generics & Collection | The creation and usage of generic classes, usage of different types of interfaces (e.g., List Interface, Deque Interface, Map Interface, and Set Interface), and comparison of objects using interfaces (e.g., java.util.Comparator, and java.lang.Comparable). |
| **[K9]** Functional Interface | The development of code that uses different versions of defined functional interfaces (e.g., primitive, binary, and unary) and user-defined functional interfaces. |
| **[K10]** Stream API | The development of code with lambda expressions and Stream APIs. This includes developing code to extract data from an object using peek() and map() methods, searching for data with search methods (e.g., findFirst) of the Stream classes, sorting a collection using Stream API, iterating code with foreach of Stream, and saving results to a collection using the collect method. |
| **[K11]** Exception | The creation of try-catch blocks, the usage of multiple catch blocks, the usage of try-with-resources statements, the invocation of methods throwing an exception, and the use of assertion for testing invariants. |
| **[K12]** Date time API | This KU refers to create and manage date-based and time-based events using Instant, Period, Duration, and TemporalUnit, and work with dates and times across timezones and manage changes resulting from daylight savings, including Format date and times values. |
| **[K13]** IO | Reading and writing data from console and files and using basic Java input-output packages (e.g., java.io.package). |
| **[K14]** NIO | Interacting files and directories with the new non-blocking input/output API (e.g., using the Path interface to operate on file and directory paths). Performing other file-related operations (e.g., read, delete, copy, move, and managing metadata of a file or directory). |
| **[K15]** String Processing | The knowledge about searching, parsing, replacing strings using regular expressions, and using string formatting. |
| **[K16]** Concurrency | The knowledge about functionalities that are related to thread execution and parallel programming. Some of these functionalities include: creating worker threads using Runnable and Callable classes, using an ExecutorService to concurrently execute tasks, using synchronized keyword and java.util.concurrent.atomic package to control the order of thread execution, using java.util.concurrent classes and using a parallel fork/join framework. |
| **[K17]** Database | The creation of database connection, submitting queries and reading results using the core JDBC API. |
| **[K18]** Localization | The knowledge about reading and setting the locale (Oracle defines a locale as "a specific geographical, political, or cultural region") by using the Locale object, and building a resource bundle for each locale, and loading a resource bundle in an application. |
| **[K19]** Java Persistence | The usage of object/relational mapping facilities for managing relational data in Java applications. With this knowledge, developers can learn how to map, store, update and retrieve data from relational databases to Java objects and vice versa. |
| **[K20]** Enterprise Java Bean | The knowledge about managing server-side components that encapsulate the business logic of an application. |
| **[K21]** Java Message Service API | The knowledge of how to create, send, receive and read messages using reliable, asynchronous, and loosely coupled communication. |
| **[K22]** SOAP Web Service | Creating and using the Simple Object Access Protocol for sending and receiving requests and responses across the Internet using JAX-WS and JAXB APIs. |
| **[K23]** Servlet | Handling HTTP requests, parameters, and cookies and how to process them on the server sites with appropriate responses. |
| **[K24]** Java REST API | Creating web services and clients according to the Representational State Transfer architectural pattern using JAX-RS APIs. |
| **[K25]** Websocket | Creating and handling bi-directional, full-duplex, and real-time communication between the server and the web browser. |
| **[K26]** Java Server Faces | The knowledge about how to build UI component-based and event-oriented web interfaces using the standard JavaServer Faces (JSF) APIs. |
| **[K27]** Contexts and Dependency Injection (CDI) | Managing the lifecycle of stateful components using domain-specific lifecycle contexts and type-safely inject components (services) into client objects. |
| **[K28]** Batch Processing | Creating and managing long-running jobs on schedule or on demand for performing on bulk-data, and without manual intervention. |

**Fig. 2** An overview of our data collection process.

We count the occurrences of KUs for every Java source file. For example, one of the capabilities associated with the *Generics and Collections KU* is to create a generic class. Hence, if we find a declaration of a `generic` class in a given Java file `x.java` belonging to a release $R$, we increment the counter for *Generic and Collection KU* that is associated with file `X.java` by one.

## 3 Data Collection

Our data collection process is straightforward and contains four steps (Figure 2). We collect the defect dataset from the replication package of Yatish et al. [56] (Step 1) and filter-in traditional metrics (Step 2). Then, we download the source code of our studied releases (Step 3) and compute the per-file KU metrics using our custom KU detector (Step 4). In the following, we describe each step of our data collection process in more detail.

• **(Step 1) Download the defect dataset.** We choose the defect dataset that was curated by Yatish et al. [56]. Instead of relying on heuristics to identify post-release defects, the authors employ an approach that "leverages the earliest affected release that is realistically estimated by a software development team for a given defect". We download the defect dataset from the article's replication package[3]. The dataset contains post-release defect data for

---

[3] `https://github.com/awsm-research/replication-icse2019`

28 releases that span a total of 8 Java project. Table 2 presents a statistical summary of the dataset (adapted from original paper).

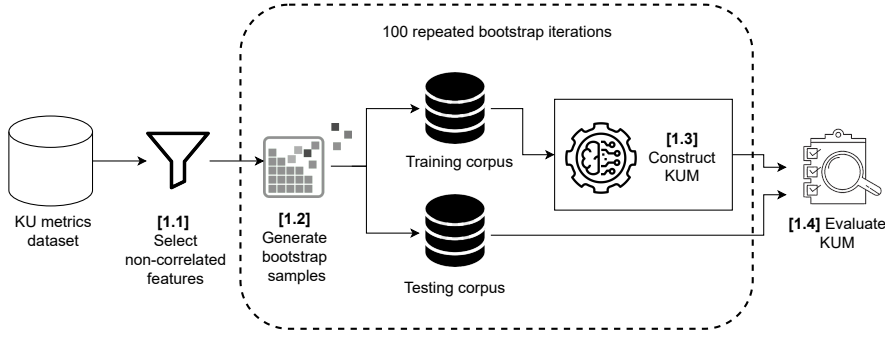**Table 2** A statistical summary of the dataset.

| Project Name | Studied Releases | Number of Java Files | Line of Code | Total Defects Reported | Rate of Defects (across releases) |
|---|---|---|---|---|---|
| ActiveMQ | 5.0.0, 5.1.0, 5.2.0, 5.3.0 | 1,884-3,420 | 142-299K | 3,157 | 6%-15% |
| Derby | 10.2.1.6, 10.3.1.4, 10.5.1.1 | 1,963-2,705 | 412-533K | 3,731 | 14%-33% |
| Groovy | 1.5.7, 1.6.0.Beta_1, 1.6.0.Beta_2 | 755-884 | 74-90K | 3,943 | 3%-8% |
| HBase | 0.94.0, 0.95.0, 0.95.2 | 1,059-1,834 | 246-534K | 5,360 | 20%-26% |
| Hive | 0.9.0, 0.10.0, 0.12.0 | 1,416-2,662 | 287-563K | 3,306 | 8%-19% |
| JRuby | 1.1, 1.4.0, 1.5.0, 1.7.0 | 731-1,614 | 105-238K | 5,475 | 5%-18% |
| Lucene | 2.3.0, 2.9.0, 3.0.0, 3.1.0 | 8,052-2,806 | 101-242K | 2,316 | 3%-24% |
| Wicket | 1.3.0.beta1, 1.3.0beta2, 1.5.3 | 1,672-2,578 | 109-165K | 3,327 | 4%-7% |

- **(Step 2) Filter-in traditional metrics (product, process, and ownership)** We filter-in the traditional metrics from the defect dataset. This filtered dataset contains 54 product, 5 process, and 6 code ownership metrics (65 total) calculated for each Java source code file. Product metrics were calculated using the Understand tool [52]. Examples include the number of lines of code (CountLineCode), number of declared public methods in a class (CountDeclMethodPublic), and number of other classes to which a class is coupled to (CountClassCoupled). The process metrics are: the number of commits (COMM), the number of lines added (ADDED_LINES), the number of lines deleted (DEL_LINES), the number of active developers (ADDEV) and the number of distinct developers (DDEV). The ownership metrics are: number of minor authors (developers contributing less than 5% of the code in a file), number of major authors (developers contributing 5% or more of the code in a module), and ownership ratio (the proportion of lines of code written by the developer who has the highest contribution of lines of code on the module). These ownership metrics are calculated for two granularities of code changes: lines of code (MINOR_LINE, MAJOR_LINE, and OWN_LINE) and commit (MINOR_COMMIT, MAJOR_COMMIT, and OWN_COMMIT).

We refer to these 65 metrics as *traditional* because they have frequently been used in the literature to build and analyze defect prediction models [15, 24, 29, 42, 58]. The full list of the traditional metrics can be seen in B.

- **(Step 3) Download the source code of the studied releases.** To download the source code of the studied software projects, we clone the repository of each project using the command git clone `project_url` where `project_url` is the url of the repository of the project in GitHub. Then, we extract the source code of a particular release using the command `git checkout tags/tag_name`. Here, `tag_name` is the name of one of the studied release for the project.

- **(Step 4) Compute KU metrics.** We count the occurrences of each KU in each Java file from every studied release in the defect dataset using our custom KU detector (refer to Section 2.3 for details).

**Fig. 3** An overview of our approach for building and evaluating KUM.

The outputs of the data collection process are the *KU metrics dataset* and the *traditional metrics dataset* (Figure 2). These datasets are employed to build several defect prediction models as part of our RQs.

## 4 Findings

In this section, we address our research questions. For each question, we discuss our motivation for studying it, the approach that we used to answer it, and our findings.

### 4.1 RQ1: How effective are KU features for predicting post-release defects?

**Motivation.** Traditionally, defect prediction has relied heavily on process, product, and code ownership metrics [20, 38, 53, 59] (traditional metrics). KUs offer a new lens through which the code of software systems can be analyzed, since they capture system traits that are not captured by those traditional metrics. More generally, we believe that exploring novel types of data for feature engineering is as important as developing or assessing new machine learning approaches for defect prediction. The more indicators researchers have, the better equipped they are to understand defects and hypotheisze their underlying causes.

**Approach.** First, we build a defect prediction model using the KU metrics dataset for each of the studied releases. We treat each KU metric as a model feature and we refer to the entire set of KU features as the *KU feature set*. Also, we refer to this KU model as KUM . Next, we evaluate the performance of KUM based on an *out-of-sample bootstrap* model validation technique. Finally, we compare the performance of KUM to baseline models that are built using the traditional metrics dataset. In the following, we explain our approach in more detail:

– *(Step 1) Construct KUM.* In this step, we construct and evaluate KUM. KUM is a binary classification model that outputs either *clean* or *defective* for an input Java file. An overview of our approach for building KUM can be seen in Figure 3.

– *(Step 1.1) Select non-correlated features.* Prior to building any defect prediction model, it is recommended to mitigate correlated features [23] to better interpret the results of a model. Therefore, we want to check the performance of our studied models using non-correlated features. We select non-correlated features using AutoSpearman [22]. AutoSpearman automatically selects non-correlated features based on two analyses: (1) a Spearman rank correlation test and (2) a Variance Inflation Factor (VIF) analysis. Applying AutoSpearman[4] to the KU feature set resulted in the selection of all features, as no feature correlations were detected.

– *(Step 1.2) Generate bootstrap samples.* To estimate the performance of our models in practice, we employ the *out-of-sample bootstrap* validation technique with 100 repetitions [12, 55]. In this technique, we randomly sample dataset with replacement to create 100 bootstrap samples. Each bootstrap sample has the same number of data points as the original dataset, but some points may appear multiple times while others may be excluded. These excluded data points from each bootstrap sample form the corresponding *out-of-sample data*, which we use for testing our models. We train our models on each bootstrap sample and test its performance on the associated out-of-sample data.

– *(Step 1.3) Construct KUM.* We use the KU feature set to construct KUM . Similarly to Yatish et al. [56], we use a random forest classifier with default parameters values (`scikit-learn` python implementation). For a detailed analysis of the influence of classifier choice and hyperparameter tuning on model performance, please refer to C.

– *(Step 1.4) Evaluate KUM.* To evaluate the accuracy of our defect prediction models, we calculate the *Area Under the receiver operator characteristic Curve* (AUC). AUC is a robust, threshold-independent measure [54]. AUC value ranges from 0 to 1, with 0.5 being as good as random guessing. The AUC values closer to one indicate a model that is good at distinguishing the class of interest (defective) from the other class (clean). Values closer to zero indicate a model that is performing worse than random guessing, consistently misclassifying instances by predicting the opposite class (i.e., labeling defective instances as clean and vice versa). Since we are employing an *out-of-sample bootstrap* model validation with 100 repetitions, we obtain 100 AUC values for each prediction model that we evaluate.

– *(Step 2) Construct baseline models.* Our main goal in RQ1 is to understand the efficacy of KUs in predicting software defects. To construct baseline models, we use the traditional metrics dataset that we extracted from our data collection. We construct four baseline model: (i) a model built with product metrics (PROD), (ii) a model built with process metrics (PROC), (iii) a model

---

[4]  https://xai4se.github.io/defect-prediction/data-preprocessing.html

built with code ownership metrics (OWN) and (iv) a model built with product, process, and ownership metrics (i.e., all the traditional metrics – TM). These baseline models are constructed following the same procedure described in Step 1.

− *(Step-3) Compare models' performance.* To check whether the performance difference between a given pair of models (e.g., KUM vs PROD) is statistically significant, we apply the Wilcoxon signed-rank test. The Wilcoxon signed-rank test is a non-parametric statistical test for paired data. We infer that there is a statistically significant difference between the two input distributions when the p-value computed by the test is smaller than 0.05 (i.e., $alpha = 0.05$). In addition, we also use the Cliff's delta ($d$) effect size measure to quantify the practical difference between the two distributions [33]. We use the following thresholds for interpreting $d$ [50]: *negligible* for $|\delta| \leq 0.147$, *small* for $0.147 < |\delta| \leq 0.33$, *medium* for $0.33 < |\delta| \leq 0.474$, and *large* otherwise.

To evaluate a model's performance improvement compared to a baseline, we define a metric called *normalized AUC improvement*. This metric measures the extent of improvement achieved by the model relative to the maximum possible improvement. It is calculated using the following formula:

$$Normalized\ AUC\ improvement = \frac{AUC\ of\ Proposed Model - AUC\ of\ Baseline Model}{1 - AUC\ of\ Baseline Model} \times 100\%$$

Here, 1 represents the highest possible AUC score. For example, to calculate the normalized AUC improvement of a model M over baseline B, where the AUC of M is 0.81 and the AUC of B is 0.75, we would compute (0.81-0.75) / (1 - 0.75) × 100 % = 24%. This result indicates that M's improvement is 24% of the maximum possible AUC improvement. Since we build 100 models for each studied system (i.e., one for each bootstrap sample), we report the average normalized AUC improvement when comparing two models.

We also rank the studied prediction models (KUM , TM , PROD, PROC, OWN) using Scott-Knott ESD (SK-ESD) [18, 36, 55]. The SK-ESD is a method of multiple comparison that leverages a hierarchical clustering to partition the set of treatment values (e.g., medians or means) into statistically distinct groups with non-negligible difference [36]. We use Tim Menzies' implementation of the SK-ESD technique, which employs non-parametric statistical tests and the Cliff's Delta measure of effect size [35].

**Findings.** **Observation 1)** ***The minimum AUC of KUM is above 0.70 (i.e., AUC > 0.70), indicating that KUs are a good candidate to predict post-release defects.*** Table 3 shows the effect size and median AUC value of the prediction models for each studied release. The KUM column refers to the AUC of the models built with the KU feature set. We observe that the AUCs of KUM are always above 0.70 (minimum is 0.71). In particular, the median AUC of KUM across all the studied releases is 0.82. A defect prediction model is typically considered suitable to be used in practice when its AUC is higher than or equal to 0.70 [24]. This result highlights that our KU features are good candidates to predict post-release defects.
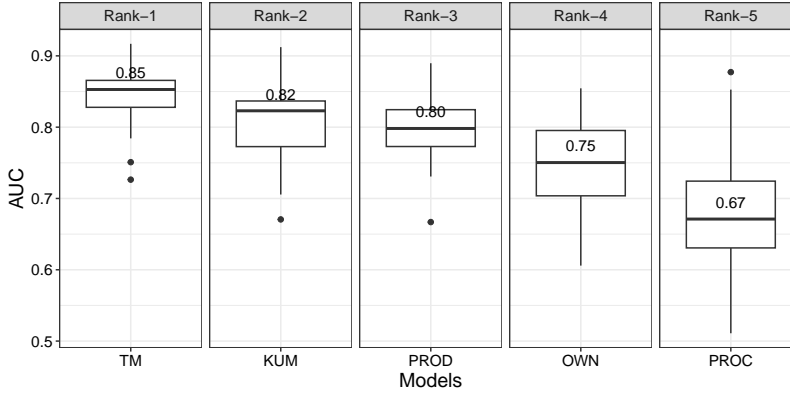
**Table 3** The effect size and median AUC value of KUM , PROD, PROC, OWN and TM
defect prediction models.

| Release | Median AUC | | | | | Statistical Test (Effect Size) | | | |
|---|---|---|---|---|---|---|---|---|---|
| | KUM | PROD | PROC | OWN | TM | KUM vs PROD | KUM vs PROC | KUM vs OWN | KUM vs TM |
| activemq-5.0.0 | 0.83 | 0.80 | 0.81 | 0.84 | 0.89 | L(0.77)* | L(0.63)* | S(-0.28)* | L(-1.00)* |
| activemq-5.1.0 | 0.83 | 0.80 | 0.62 | 0.73 | 0.82 | L(0.76)* | L(1.00)* | L(0.99)* | S(0.21)* |
| activemq-5.2.0 | 0.83 | 0.80 | 0.69 | 0.82 | 0.85 | L(0.75)* | L(1.00)* | N(0.07) | L(-0.70)* |
| activemq-5.3.0 | 0.82 | 0.80 | 0.66 | 0.81 | 0.84 | L(0.64)* | L(1.00)* | S(0.29)* | L(-0.58)* |
| activemq-5.8.0 | 0.85 | 0.83 | 0.68 | 0.74 | 0.86 | L(0.61)* | L(1.00)* | L(1.00)* | M(-0.41)* |
| derby-10.2.1.6 | 0.74 | 0.77 | 0.78 | 0.79 | 0.86 | L(-0.73)* | L(-0.95)* | L(-0.99)* | L(-1.00)* |
| derby-10.3.1.4 | 0.84 | 0.82 | 0.71 | 0.75 | 0.86 | S(0.18)* | L(1.00)* | L(1.00)* | L(-0.72)* |
| derby-10.5.1.1 | 0.83 | 0.82 | 0.60 | 0.75 | 0.85 | M(0.44)* | L(1.00)* | L(1.00)* | L(-0.87)* |
| groovy-1.5.7 | 0.81 | 0.77 | 0.53 | 0.71 | 0.78 | M(0.38)* | L(0.93)* | L(0.68)* | S(0.24)* |
| groovy-1.6.BETA_1 | 0.82 | 0.80 | 0.72 | 0.75 | 0.87 | S(0.23)* | L(0.79)* | L(0.76)* | L(-0.64)* |
| groovy-1.6.BETA_2 | 0.85 | 0.89 | 0.75 | 0.77 | 0.91 | L(-0.5')* | L(0.90)* | L(0.83)* | L(-0.64)* |
| hbase-0.94.0 | 0.71 | 0.78 | 0.64 | 0.76 | 0.83 | L(-0.9')* | L(0.97)* | L(-0.78)* | L(-1.00)* |
| hbase-0.95.0 | 0.71 | 0.69 | 0.67 | 0.70 | 0.82 | L(0.90)* | L(0.82)* | S(0.22)* | L(-1.00)* |
| hbase-0.95.2 | 0.71 | 0.65 | 0.63 | 0.69 | 0.73 | L(0.69)* | L(0.83)* | M(-0.39)* | L(-0.95)* |
| hive-0.10.0 | 0.72 | 0.78 | 0.67 | 0.79 | 0.85 | L(-0.90)* | M(0.43)* | L(-0.93)* | L(-1.00)* |
| hive-0.12.0 | 0.72 | 0.79 | 0.62 | 0.80 | 0.85 | L(-0.97)* | L(0.94)* | L(-0.98)* | L(-1.00)* |
| hive-0.9.0 | 0.79 | 0.86 | 0.57 | 0.70 | 0.91 | L(-0.92)* | L(1.00)* | L(0.98)* | L(-1.00)* |
| jruby-1.1 | 0.86 | 0.82 | 0.66 | 0.83 | 0.84 | L(0.66)* | L(1.00)* | M(0.42)* | M(0.49)* |
| jruby-1.4.0 | 0.83 | 0.82 | 0.67 | 0.73 | 0.85 | M(0.45)* | L(1.00)* | L(0.98)* | M(-0.44)* |
| jruby-1.5.0 | 0.84 | 0.82 | 0.75 | 0.82 | 0.86 | M(0.45)* | L(0.95)* | M(0.39)* | S(-0.28)* |
| jruby-1.7.0.preview1 | 0.84 | 0.81 | 0.51 | 0.76 | 0.85 | L(0.74)* | L(1.00)* | L(0.90)* | M(-0.39)* |
| lucene-2.3.0 | 0.82 | 0.84 | 0.88 | 0.85 | 0.92 | M(-0.42)* | L(-0.81)* | L(-0.74)* | L(-1.00)* |
| lucene-2.9.0 | 0.78 | 0.78 | 0.68 | 0.69 | 0.82 | M(0.41)* | L(1.00)* | L(1.00)* | L(-0.77)* |
| lucene-3.0.0 | 0.91 | 0.89 | 0.71 | 0.68 | 0.90 | L(0.81)* | L(1.00)* | L(1.00)* | M(0.46)* |
| lucene-3.1.0 | 0.73 | 0.72 | 0.64 | 0.61 | 0.75 | S(0.27)* | L(0.87)* | L(0.96)* | L(-0.56)* |
| wicket-1.3.0-beta2 | 0.82 | 0.80 | 0.74 | 0.70 | 0.83 | M(0.44)* | L(0.94)* | L(1.00)* | M(-0.43)* |
| wicket-1.3.0-incubating-beta-1 | 0.86 | 0.84 | 0.85 | 0.79 | 0.87 | M(0.13)* | M(0.38)* | L(0.62)* | L(-0.66)* |
| wicket-1.5.3 | 0.81 | 0.77 | 0.53 | 0.70 | 0.79 | L(0.73)* | L(1.00)* | L(0.97)* | M(0.39)* |
| Median | 0.82 | 0.80 | 0.67 | 0.75 | 0.85 | | | | |

1. The green cells indicate cases where KUM outperforms a baseline model with a non-negligible effect size.
2. Cliff's Delta effect size: L – Large, M – Medium, S – Small, and N – Negligible
3. The statistical tests with a p-value less than 0.05 are marked with an asterisk (*) in the effect size results.



**Fig. 4** The distribution of AUC of KUM and other studied models. The models are grouped
based on their performance rankings determined by the Scott-Knott ESD (SK-ESD) method,
where a lower SK-ESD rank indicates a better-performing model.

**Observation 2) *KUM outperforms all baseline models except for
TM.*** Figure 4 presents the SK-ESD ranks of the studied models across the
studied releases. KUM ranks second, outperforming the baseline models that
are built with one group of the traditional metrics (PROD, OWN and PROC).
The median AUC values for these baseline models are 0.80 (PROD), 0.75

(OWN), and 0.67 (PROC), whereas, the median AUC for KUM is 0.82. Indeed, analysis of Table 3 indicates that KUM consistently outperforms these three baseline models in the majority of the studied releases. For example, KUM outperforms PROD in 21 out of the 28 studied releases, typically with either a medium or large effect size (note green rows in Table 3 of the column named "KUM vs PROD"). KU metrics can thus be interpreted as a set of *product* metrics with higher predictive power than traditional product metrics. Additionally, we calculate the normalized AUC improvement of KUM over PROD, PROC, and OWN for each release. The normalized AUC improvement of KUM over PROD ranges from 3.6% to 22.2% (median 11.4%), improvement over PROC ranges from 6.7% to 68.9% (median 42.4%), and improvement over OWN ranges from 3.3% to 71.9% (median 33.3%). In summary, KUM not only outperforms a model built on the same raw data (PROD), but also those built on process (PROC) and ownership (OWN) data.

Finally, while TM does outperform KUM in 23 out of the 28 releases (Table 3), KUM achieves a median AUC of 0.82 in those 23 releases. This result indicates that KUM would still be a robust classifier for those releases in practice.

---

**Summary**

**RQ1: How effective are KU features for predicting post-release defects?**

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

KUs are good candidates for predicting defect-prone code. In particular:

• KUM can classify post-release defects with a minimum AUC of 0.71 and a median AUC of 0.82.
• KUM outperforms the baseline models that are built with each individual group of traditional metrics (PROC, PROD and OWN). The normalized AUC improvement of KUM over PROD ranges from 3.6% to 22.2% (median 11.4%), improvement over PROC ranges from 6.7% to 68.9% (median 42.4%), and improvement over OWN ranges from 3.3% to 71.9% (median 33.3%).
• While KUM does not outperform TM (model built with all traditional metrics), KUM maintains strong AUC values.

---

## 4.2 RQ2: What are the most important KU features for predicting post-release defects?

**Motivation.** Analyzing feature importances helps pinpoint the features that contribute most significantly to predicting post-release defects. This analysis not only enhances our understanding of the model's decision-making process but also highlights the practical relevance of specific KUs.
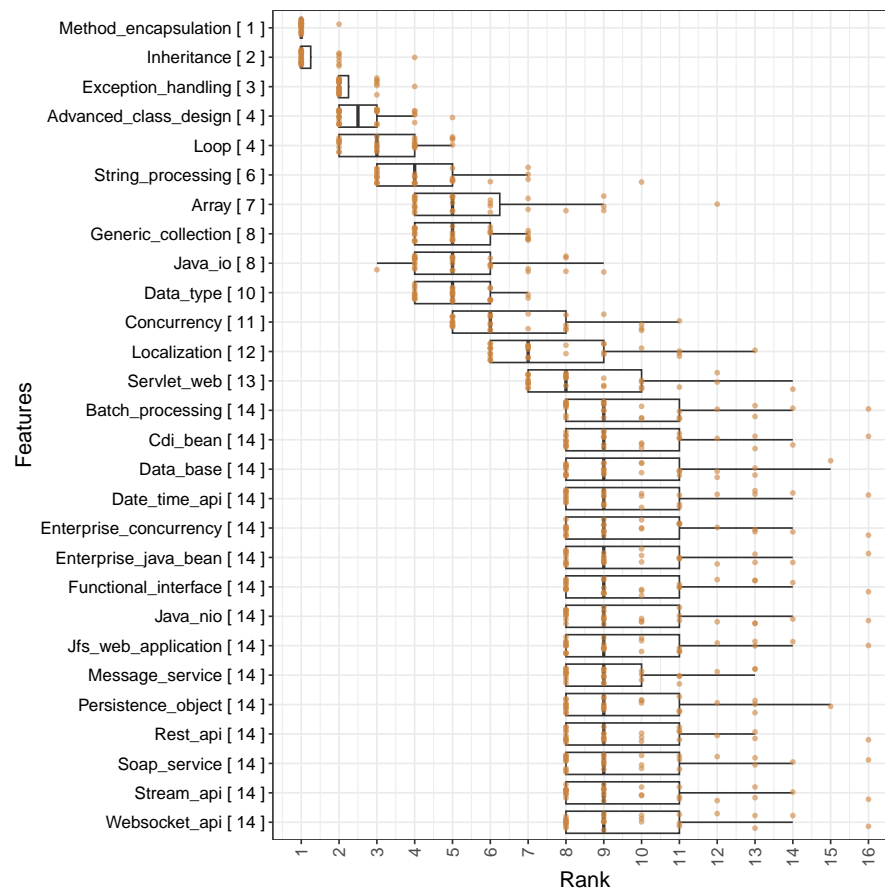
**Approach.** To determine the KU features that have the most significant impact on predictions, we conduct a model interpretation analysis utilizing SHAP (SHapley Additive exPlanation). SHAP is a widely adopted method for model interpretation that leverages game theory to provide an optimal framework for quantifying each feature's contribution to the model's predictions [37]. SHAP

values can be positive or negative, signifying the direction of the feature's influence on the prediction. A positive SHAP value indicates that the feature increases the likelihood of the predicted outcome (e.g., classifying a file as defective in a defect prediction model). In contrast, a negative SHAP value means that the feature reduces the likelihood of the predicted outcome or increases the likelihood of the opposite class (i.e., clean class). The strength of a feature is determined by the magnitude of SHAP values. Specifically, the absolute SHAP value of a feature indicates its influence on the prediction, with higher values signifying greater importance. We use the `shapper` python package to perform the SHAP analysis [9].

We compute feature importances for KUM as follows. First, we apply SHAP to KUM and calculate the absolute SHAP value for each feature across every file (i.e., record) in the dataset of a release. Next, for each feature, we sum the SHAP values of all files (i.e., all records) in the dataset. We follow this approach for all the 100 bootstraps of every release. As a result we obtain a distribution of the sum of SHAP values for every feature across all releases. Then, we apply the SK-ESD technique [36] on the distribution of the sum of SHAP values generated from the 100 bootstraps of every release. As a result, we obtain a feature ranking for each release. Since there can be ties, two features with same importance are given the same rank. By taking all releases, we can produce a feature rank distribution for each feature. To summarize the results across all releases, we apply the SK-ESD technique again on the feature rank distributions, similarly to Ghotra et al. [18]. As a result, we obtain a final rank for each feature.

**Findings. Observation 3)** *Method & Encapsulation is the most important KU feature across all studied releases of KUM for predicting post-release defects.* Figure 5 shows the rank distribution of each feature (final rank is shown in brackets). We observe that Method & Encapsulation is the top-ranked feature in KUM. Indeed, the boxplot reveals that Method & Encapsulation is the most important feature in most of our studied releases as the median value is one and the interquartile range is narrow (most ranks are one). Inheritance is ranked second and Exception Handling is the third-ranked feature. We then have two features ranked forth, namely: Advanced Class Design and Loop. Some advanced KUs derived from Java SE II certification exams also rank among the top ten most important features. For instance, Generics and Collections and Java IO are both ranked eighth. KUs related to building enterprise-level applications in Java, such as Enterprise JavaBeans and Java Message Service (JMS), which are derived from the Java EE certification exam, are ranked lower in importance (typically 14th). Upon closer inspection, we note that these advanced Java EE features are less frequently used (or even not used at all) in the studied releases. In contrast, fundamental KUs like Data Type and Loop, derived from the Java SE I exam, are more prevalent.

**Fig. 5** The rank distribution of features of KUM for defect classification. The number inside the square brackets indicates the final rank of the feature after applying Scott-Knott ESD for the second time.

**Summary**

**RQ2: What are the most important KU features for predicting post-release defects?**

Method & Encapsulation, Inheritance, Exception Handling, Advanced Class Design, and Loop are the top five most important KU features for predicting post-release defects across all studied releases.

**Fig. 6** The distribution of AUC of the combined KUM+TM and other studied models. The models are grouped based on their performance rankings determined by the Scott-Knott ESD (SK-ESD) method, where a lower SK-ESD rank indicates a better-performing model.

### 4.3 RQ3: Can KUM be made more accurate by combining its features with those from TM?

**Motivation.** This research question investigates whether the predictive performance of KUM for post-release defects can be improved by integrating them with traditional metrics. The underlying motivation is based on the hypothesis that combining diverse sets of features can provide a more comprehensive understanding of the factors influencing post-release defects. Such integration may enable the combined model to capture a wider range of predictive indicators, thereby enhancing its accuracy in identifying potential defects.

**Approach.** We build a combined model named KUM+TM for predicting post-release defects using all studied features of KUM and TM. Following Step 1 of RQ1 (Figure 3), we build KUM+TM models for all studied releases. We rank the studied prediction models (KUM, TM, PROD, PROC, OWN and KUM+TM) using SK-ESD [18, 36, 55]. We further compare the performance of KUM+TM model with the performance of KUM and TM models for each of the studied releases. We also analyze feature importances of KUM+TM model using the same SHAP analysis method explained in Section 4.2.

**Findings.   Observation 4)** *Combining KU features with traditional metrics results in a better performing model for predicting post-release defects.* Figure 6 depicts the performance of the studied models. We observe that the combined model KUM+TM is the top performing model. The median AUC of KUM+TM is 0.89, representing a normalized improvement of 36.4% over TM and 63.6% over KUM. We further examine the performance of KUM+TM for individual releases (Table 4). We observe that KUM+TM models outperform TM and KUM in all 28 studied releases, with the majority showing a large effect size (highlighted in green in the columns "KUM+TM vs

**Table 4** The effect size and median AUC value of KUM, TM, and KUM+TM defect classification models.

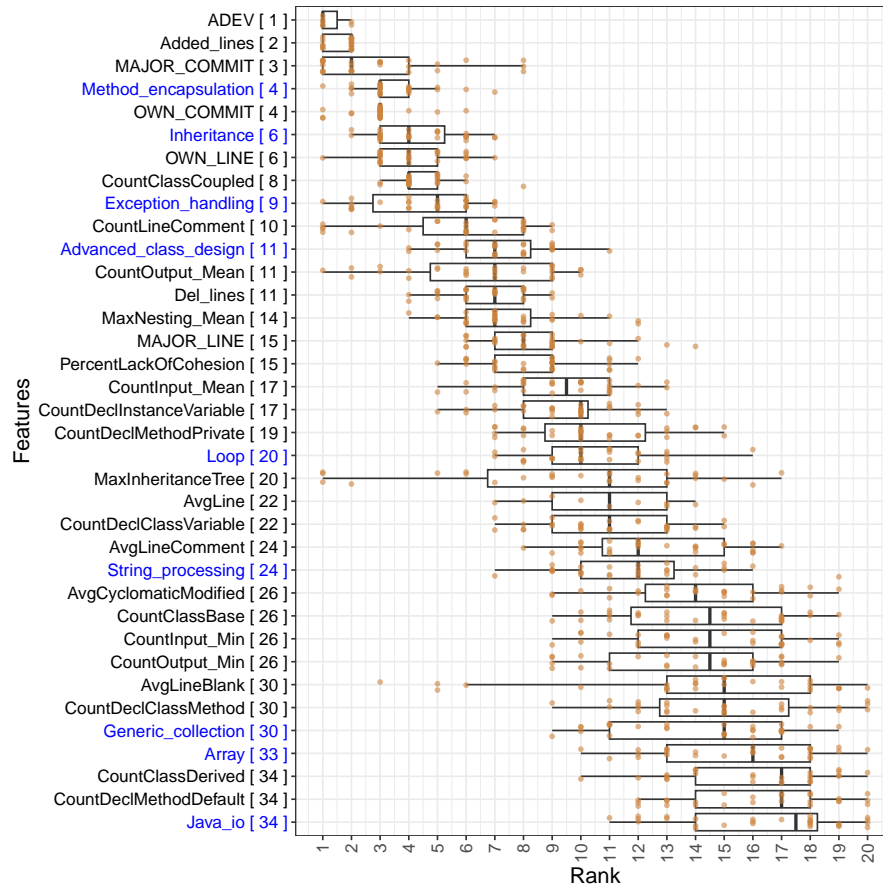| Release | Median AUC | | | Statistical Test (Effect Size) | |
|---|---|---|---|---|---|
| | KUM+TM | KUM | TM | KUM+TM vs KUM | KUM+TM vs TM |
| activemq-5.0.0 | 0.92 | 0.83 | 0.89 | L(1.00)* | L(0.94)* |
| activemq-5.1.0 | 0.86 | 0.83 | 0.82 | L(0.68)* | L(0.81)* |
| activemq-5.2.0 | 0.89 | 0.83 | 0.85 | L(0.98)* | L(0.85)* |
| activemq-5.3.0 | 0.87 | 0.82 | 0.84 | L(0.97)* | L(0.81)* |
| activemq-5.8.0 | 0.89 | 0.85 | 0.86 | L(0.97)* | L(0.93)* |
| derby-10.2.1.6 | 0.89 | 0.74 | 0.86 | L(1.00)* | L(0.95)* |
| derby-10.3.1.4 | 0.90 | 0.84 | 0.86 | L(1.00)* | L(0.99)* |
| derby-10.5.1.1 | 0.89 | 0.83 | 0.85 | L(1.00)* | L(0.95)* |
| groovy-1.5.7 | 0.83 | 0.81 | 0.78 | S(0.18)* | M(0.40)* |
| groovy-1.6.BETA_1 | 0.90 | 0.82 | 0.87 | L(0.81)* | M(0.42)* |
| groovy-1.6.BETA_2 | 0.95 | 0.85 | 0.91 | L(0.96)* | L(0.74)* |
| hbase-0.94.0 | 0.86 | 0.71 | 0.83 | L(1.00)* | L(0.60)* |
| hbase-0.95.0 | 0.84 | 0.71 | 0.82 | L(1.00)* | L(0.68)* |
| hbase-0.95.2 | 0.76 | 0.67 | 0.73 | L(1.00)* | L(0.79)* |
| hive-0.10.0 | 0.88 | 0.72 | 0.85 | L(1.00)* | L(0.87)* |
| hive-0.12.0 | 0.88 | 0.72 | 0.85 | L(1.00)* | L(0.80)* |
| hive-0.9.0 | 0.94 | 0.79 | 0.91 | L(1.00)* | L(0.96)* |
| jruby-1.1 | 0.90 | 0.86 | 0.87 | L(0.58)* | L(0.63)* |
| jruby-1.4.0 | 0.88 | 0.83 | 0.85 | L(0.84)* | L(0.67)* |
| jruby-1.5.0 | 0.88 | 0.84 | 0.86 | L(0.68)* | L(0.54)* |
| jruby-1.7.0.preview1 | 0.89 | 0.84 | 0.85 | L(0.93)* | L(0.83)* |
| lucene-2.3.0 | 0.95 | 0.82 | 0.92 | L(1.00)* | L(0.91)* |
| lucene-2.9.0 | 0.85 | 0.78 | 0.82 | L(0.99)* | L(0.88)* |
| lucene-3.0.0 | 0.93 | 0.91 | 0.90 | L(0.67)* | L(0.91)* |
| lucene-3.1.0 | 0.80 | 0.73 | 0.75 | L(0.96)* | L(0.74)* |
| wicket-1.3.0-beta2 | 0.87 | 0.82 | 0.83 | L(0.93)* | L(0.81)* |
| wicket-1.3.0-incubating-beta-1 | 0.90 | 0.86 | 0.87 | L(0.95)* | L(0.80)* |
| wicket-1.5.3 | 0.83 | 0.81 | 0.79 | M(0.41)* | L(0.72)* |
| Median | 0.89 | 0.82 | 0.85 | | |

1. The green color rows show cases where KUM+TM significantly outperform the studied models with a non-negligible effect size.
2. Cliff's Delta effect size: L – Large, M – Medium, S – Small, and N – Negligible
3. The statistical tests with a p-value less than 0.05 are marked with an asterisk (*) in the effect size results.

KUM" and "KUM+TM vs TM"). We also calculate the normalized AUC improvement for each of the studied releases. The normalized AUC improvement of KUM+TM over TM ranges from 11.1% to 44.4% (median 22.5%), while the improvement over KUM ranges from 10.5% to 72.2% (median 31.5%). We thus conclude that KUM+TM demonstrates superior performance for predicting post-release defects, making it more effective than either of the individual models.

**Observation 5)** *ADEV is the top ranked feature for predicting post-release defects. Method & Encapsulation is the highest ranked KU feature (fourth).* Figure 7 presents the rank distribution of each feature (final rank is shown in brackets) of KUM+TM across the studied releases. The top two most important features are process metrics: ADEV (i.e., the number of active developers) and Added_Lines (the normalized count of lines added to the file). Major Commit, an ownership metric, is ranked third. We then have two features that are ranked as the forth most important feature. One of these features is Method & Encapsulation, which is a KU feature. The other KU features in the top ten set are Inheritance (sixth) and Exception Handling (ninth).
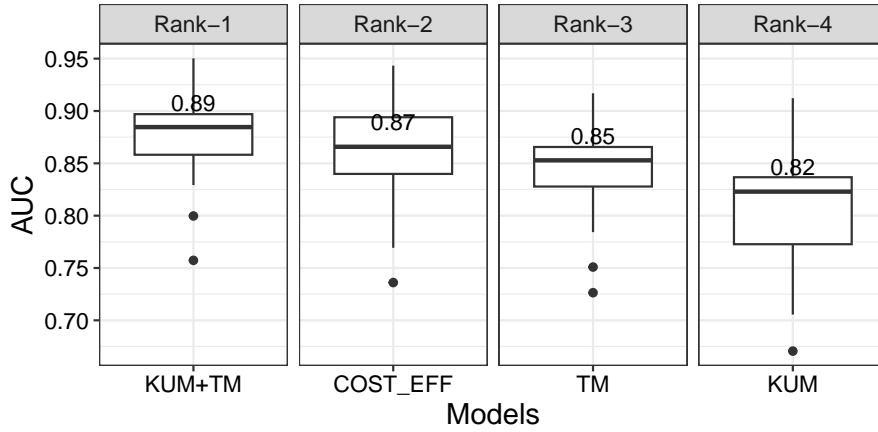
**Fig. 7** The rank distribution of features of KUM+TM. The number inside the square brackets indicates the final rank of the feature after applying Scott-Knott ESD for the second time. The KUs are highlighted with the blue colored font and traditional metrics are highlighted with black colored font.

**Summary**

**RQ3: Can KUM be made more accurate by combining its features with those from TM?**

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

Yes. KUM+TM model outperforms both TM and KUM with a non-negligible effect size in every studied system. In particular:

- KUM+TM achieves an impressive median AUC of 0.89.
- The normalized AUC improvement of KUM+TM over TM ranges from 11.1% to 44.4% (median 22.5%) and the improvement over KUM ranges from 10.5% to 72.2% (median 31.5%) across the studied releases.
- In KUM+TM, the the number of active developers (ADEV) is the most important feature. Method & Encapsulation is the top-ranked KU feature (fourth).

**Fig. 8** The distribution of AUC of COST_EFF, KUM, TM, and KUM+TM across the studied releases. The models are grouped based on their performance rankings determined by the Scott-Knott ESD (SK-ESD) method, where a lower SK-ESD rank indicates a better-performing model.

## 4.4 RQ4: Can a cost-effective model be created by combining specific features from KUM and TM?

**Motivation.** Building predictive models with a broad set of features often yields high performance, but this approach can be resource-intensive, requiring significant efforts in data collection and processing. These challenges highlight the importance of developing cost-effective models that achieve strong predictive performance while relying on a minimal, carefully selected set of features. Such models not only reduce resource requirements but also enhance practicality and scalability in real-world.

**Approach.** We focus on a subset of KU metrics and traditional metrics to build a cost-effective model. Specifically, we select the top five KU features (out of 28) from KUM and the top five metrics (out of 65) from TM. This results in the construction of our cost-effective model (COST_EFF) using only 10 features (a reduction of approximately 89% in the number of features compared to KUM+TM). The selected features are as follows: Method & Encapsulation KU, Exception KU, Advanced Class Design KU, Inheritance, Loop KU, ADEV, Added Lines, Major Commit, CountClassCoupled and CountLineComment.

To build COST_EFF models, we follow the approach for model construction presented in Section 4.1 (Figure 3). We then evaluate COST_EFF and compare its performance to that of KUM, TM and KUM+TM. The comparison involves ranking the models according to the same SK-ESD method detailed in Section 4.2.

**Table 5** The effect size and median AUC value of COST_EFF, TM, KUM and KUM+TM defect prediction models.

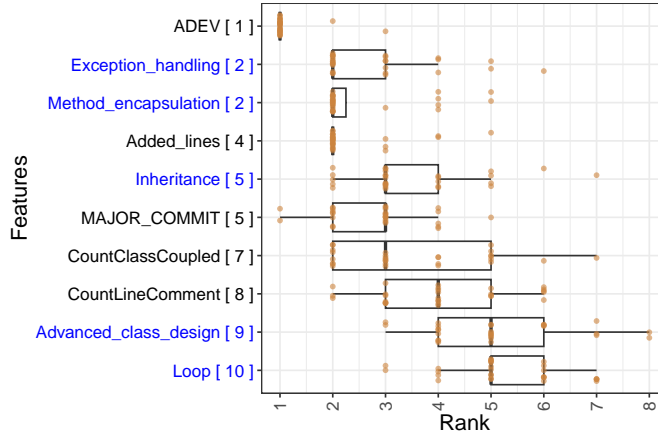| Release | Median AUC | | | | Statistical Test (Effect Size) | | |
|---|---|---|---|---|---|---|---|
| | COST_EFF | TM | KUM | KUM+TM | COST_EFF vs TM | COST_EFF vs KUM | COST_EFF vs KUM+TM |
| activemq-5.0.0 | 0.91 | 0.89 | 0.83 | 0.92 | L(0.77)* | L(1.00)* | L(-0.58)* |
| activemq-5.1.0 | 0.84 | 0.82 | 0.83 | 0.86 | L(0.53)* | M(0.33)* | M(-0.43)* |
| activemq-5.2.0 | 0.88 | 0.85 | 0.83 | 0.89 | L(0.73)* | L(0.97)* | S(-0.25)* |
| activemq-5.3.0 | 0.86 | 0.84 | 0.82 | 0.87 | L(0.52)* | L(0.88)* | L(-0.49)* |
| activemq-5.8.0 | 0.87 | 0.86 | 0.85 | 0.89 | M(0.43)* | L(0.72)* | L(-0.81)* |
| derby-10.2.1.6 | 0.87 | 0.86 | 0.74 | 0.89 | M(0.41)* | L(1.00)* | L(-0.84)* |
| derby-10.3.1.4 | 0.83 | 0.86 | 0.84 | 0.90 | L(-0.84)* | S(-0.23)* | L(-1.00)* |
| derby-10.5.1.1 | 0.85 | 0.85 | 0.83 | 0.89 | N(0.07)* | L(0.66)* | L(-0.98)* |
| groovy-1.5.7 | 0.79 | 0.78 | 0.81 | 0.83 | S(0.17)* | S(-0.16)* | S(-0.31)* |
| groovy-1.6.BETA_1 | 0.90 | 0.87 | 0.82 | 0.90 | L(0.52)* | L(0.87)* | N(0.06) |
| groovy-1.6.BETA_2 | 0.92 | 0.91 | 0.85 | 0.95 | S(0.22)* | L(0.74)* | L(-0.59)* |
| hbase-0.94.0 | 0.86 | 0.83 | 0.71 | 0.86 | L(0.78)* | L(1.00)* | S(0.22)* |
| hbase-0.95.0 | 0.84 | 0.82 | 0.71 | 0.84 | L(0.75)* | L(1.00)* | N(0.08) |
| hbase-0.95.2 | 0.74 | 0.73 | 0.71 | 0.76 | S(0.33)* | L(0.99)* | L(-0.62)* |
| hive-0.10.0 | 0.88 | 0.85 | 0.72 | 0.88 | L(0.89)* | L(1.00)* | N(0.05) |
| hive-0.12.0 | 0.88 | 0.85 | 0.72 | 0.88 | L(0.87)* | L(1.00)* | N(0.14)* |
| hive-0.9.0 | 0.92 | 0.91 | 0.79 | 0.94 | L(0.51)* | L(1.00)* | L(-0.73)* |
| jruby-1.1 | 0.89 | 0.87 | 0.86 | 0.90 | L(0.55)* | L(0.49)* | N(-0.13)* |
| jruby-1.4.0 | 0.85 | 0.85 | 0.83 | 0.88 | N(-0.11)* | M(0.38)* | L(-0.74)* |
| jruby-1.5.0 | 0.88 | 0.86 | 0.84 | 0.88 | L(0.58)* | L(0.72)* | N(0.04) |
| jruby-1.7.0.preview1 | 0.86 | 0.85 | 0.84 | 0.89 | S(0.29)* | M(0.46)* | L(-0.77)* |
| lucene-2.3.0 | 0.94 | 0.92 | 0.82 | 0.95 | L(0.85)* | L(1.00)* | S(-0.32)* |
| lucene-2.9.0 | 0.83 | 0.82 | 0.78 | 0.85 | M(0.38)* | L(0.89)* | L(-0.69)* |
| lucene-3.0.0 | 0.90 | 0.90 | 0.91 | 0.93 | N(-0.02) | M(-0.47)* | L(-0.92)* |
| lucene-3.1.0 | 0.77 | 0.75 | 0.72 | 0.80 | S(0.25)* | L(0.72)* | L(-0.61)* |
| wicket-1.3.0-beta2 | 0.83 | 0.83 | 0.82 | 0.87 | N(-0.10)* | S(0.30)* | L(-0.81)* |
| wicket-1.3.0-incubating-beta-1 | 0.90 | 0.87 | 0.84 | 0.90 | L(0.72)* | L(0.93)* | S(-0.16)* |
| wicket-1.5.3 | 0.81 | 0.79 | 0.81 | 0.83 | M(0.40)* | N(-0.04) | L(-0.48)* |
| Median | 0.87 | 0.85 | 0.82 | 0.89 | | | |

1. The green color rows show cases where COST_EFF significantly outperform the studied models with a non-negligible effect size.
2. Cliff's Delta effect size: L – Large, M – Medium, S – Small, and N – Negligible.
3. The statistical tests with a p-value less than 0.05 are marked with an asterisk (*) in the effect size results.

**Findings.   Observation 6)** *COST_EFF outperforms both TM and KUM, achieving a median AUC of 0.87* . Figure 8 presents the AUC distribution of COST_EFF, KUM, TM, and KUM+TM. COST_EFF achieves a median AUC of 0.87, ranking second overall across the studied releases (i.e., just 0.02 absolute points behind the KUM+TM model). We further analyze the model's performance for each of the studied releases. Table 5 presents the effect size of the classification models for each studied release. COST_EFF outperforms the TM in most of the studied releases (23 out of 28), with a majority of these cases showing a large effect size. In another four releases, the two models perform the same (negligible effect size). Similarly, COST_EFF outperforms KUM in most releases. The normalized AUC improvement of COST_EFF over TM ranges from 3.7% to 31.3% (median 11.1%), while the improvement over KUM ranges from 5.6% to 66.7% (median 26.8%).

**Observation 7)** *Exception Handling and Method & Encapsulation stand out as important features in COST_EFF.* Figure 9 presents the rank distribution of each feature (final rank is shown in brackets) of COST_EFF across the studied releases. Similar to the feature importance results of KUM+TM, the top most important feature of COST_EFF is ADEV (the number of active developers). In particular, ADEV is consistently the top feature in almost all releases (note how the Q1, median, and Q3 all correspond to one). The KU features Exception Handling and Method & Encapsulation are both ranked second. The narrow interquartile range of the Method &

**Fig. 9** The rank distribution of features of COST_EFF. The number inside the square brackets indicates the final rank of the feature after applying Scott-Knott ESD for the second time. The KUs are highlighted with the blue colored font and traditional metrics are highlighted with black colored font.

Encapsulation feature, compared to the wider range of Exception Handling and other lower-ranked features, indicates that Method & Encapsulation has a consistent ranking across the studied releases. Since three KU features are included in the top-5 ranked features, we conclude that KU features have a significant contribution in the cost-effective model.
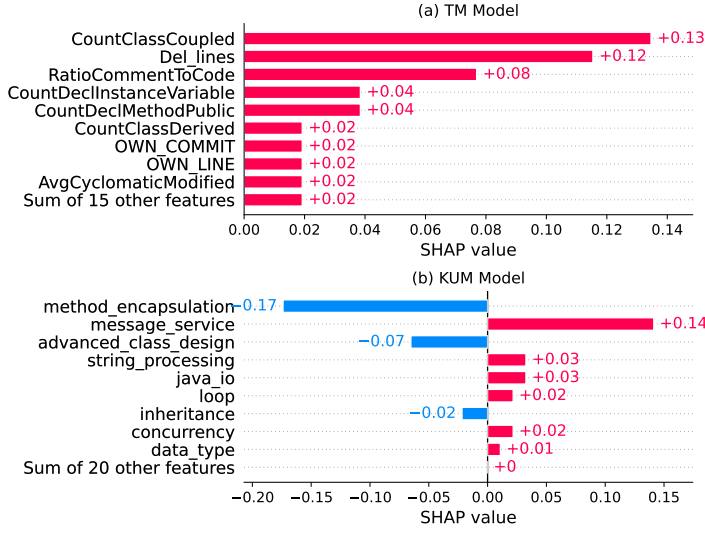
**Summary**

**RQ4: Can a cost-effective model be created by combining specific features from KUM and TM?**

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

The cost-effective model, COST_EFF, outperforms TM (as well as KUM). In particular,

• The median AUC of COST_EFF is 0.87.
• The normalized AUC improvement of COST_EFF over TM ranges from 3.7% to 31.3% (median 11.1%).
• Similar to KUM+TM, ADEV is ranked as the most important feature. Method & Encapsulation and Exception Handling are both ranked second, emphasizing the significant contribution of KU features to the performance of COST_EFF.

## 5 Discussion

In this section, we discuss the trade-offs between traditional metrics and KU metrics for predicting post-release defects and outline future directions for expanding the concept of KUs to include domain-specific and technical knowledge.

**(a) TM Model**

CountClassCoupled +0.13
Del_lines +0.12
RatioCommentToCode +0.08
CountDeclInstanceVariable +0.04
CountDeclMethodPublic +0.04
CountClassDerived +0.02
OWN_COMMIT +0.02
OWN_LINE +0.02
AvgCyclomaticModified +0.02
Sum of 15 other features +0.02

SHAP value

**(b) KUM Model**

method_encapsulation −0.17
message_service +0.14
advanced_class_design −0.07
string_processing +0.03
java_io +0.03
loop +0.02
inheritance −0.02
concurrency +0.02
data_type +0.01
Sum of 20 other features +0

SHAP value

**Fig. 10** The SHAP values of the top ten features of the TM and KUM models for their contribution to classifying the ActiveMQSession.java file as defect-prone.

## 5.1 On the trade-offs between traditional metrics and KU metrics for predicting post-release defects

Based on our findings, it is clear that KUs are **excellent** product metrics - even better than traditional product metrics (KUM outperforms PROD), which has been a long-standing standard in software engineering field. Furthermore, combining KUs with traditional metrics of TM leads to a top-performing model (KUM+TM) that achieves a median AUC of 0.89 and surpasses the already high median AUC of 0.85 achieved by TM. These findings demonstrate that KUs enhance the predictive power of traditional defect classification models.

On the other hand, we acknowledge that the traditional metrics combo (process, product and ownership metrics) remains overall stronger than just KUs (i.e., TM outperforms KUM). Feature importance analysis of the combined model KUM+TM also reveals that traditional metrics still prevail in terms of feature importances.

This duality raises the question of which set of metrics one should use in practice. To answer this question, we would like to point out that actionability is as important as predictive power in several scenarios. Since KUs offer a new to interpret post-release defects (due to its unique capability of capturing system traits), we believe that KUs have the potential to offer unique actionability insights. As an illustrative example, consider the SHAP analysis of the ActiveMQSession.java file of the activemq-5.8.0 release for TM and KUM (Figure 10). The SHAP for KUM indicates that Java Message Service API (message_service) is a relevant and important KU that positively impacts

the file classification as defective. By manually investigating the fixing patch[5], we note that one of the changed methods (run(), line 860) does use Message Service APIs (an instantiation of an object of `MessageAck` class). In addition, the defect-prone class file instantiates and uses a set of methods and subclasses of several Message Service APIs (e.g., `Session`, `Message`, `QueueSender` and `QueueReceiver`). Therefore, it becomes clear that the defect-prone class (and the defect-prone method) do operate in a context where the concept of message service matters. In fact, the associated issue AMQ-4634[6] also mentions that the problem might be related to how ActiveMQ handles transactions in a Java Message Service (JMS) environment when working with distributed transactions. In contrast, the traditional metrics of TM (e.g., CountClassCoupled) do not clearly highlight the deeper insight about the context of the defect (see the SHAP for TM in Figure 10(a)). Hence, we conclude that, in this case, KU features could provide further insights into the context in which defects happen.

In summary, in the context where predictive power is the key factor, developers should opt for COST_EFF, TM, or KUM+TM model, as those offer superior performance overall. However, if actionability and model interpretability are more important, developers may benefit from using KUM. Further research should be carried out to better understand the trade-offs between using traditional metrics and KUs (e.g., by investigating additional subject systems where Java EE KUs are more prevalent).

## 5.2 The road ahead

**Application domains.** In this paper, we leverage KUs of the Java programming language to predict post-release defects. Our results clearly demonstrate that KUs are good features for defect prediction. Our prior work also showed that KUs can effectively enhance reviewer recommendation systems [2], further highlighting their versatility and potential in software engineering tasks. While our studies lay out a strong foundation for the use of KUs in software analysis, we see these contributions as the beginning of a broader research journey rather than an endpoint. A promising direction for future work would be to explore the application of KUs in supporting line-level defect prediction [17]. For instance, KUs could be mapped to line-level code features, such as the presence of specific constructs or the usage of critical APIs. This adaptation could provide a more granular understanding of defect-prone areas and further extend the applicability of KUs in fine-grained software system analysis.

**Alternative KU elicitation methods.** Our current approach to eliciting KUs and the associated feature engineering resulted in high-performing models for predicting post-release defects. However, the process of manually extracting KUs from certification exams is labor-intensive and does not easily scale

---

[5] `https://shorturl.at/Kkc3L`

[6] `https://issues.apache.org/jira/browse/AMQ-4634`

to other programming languages or contexts. To address this limitation, researchers should explore more efficient and scalable methods for eliciting KUs. One promising avenue is leveraging Large Language Models (LLMs) to automate the elicitation and summarization of KUs [41]. LLMs could be trained or fine-tuned to extract KUs not only from certification exams but also from a broader range of resources, such as classic textbooks, programming tutorials, and domain-specific documentation. This approach could significantly expand the applicability of KUs across different languages and platforms. It is also worth noting that certification exams do exist for various programming languages and could serve as a foundation for eliciting KUs in those contexts. However, automating and generalizing the elicitation process by incorporating diverse sources would make the process more practical and adaptable for a wider range of programming languages and scenarios.

**Beyond KUs of programming languages.** Researchers should consider expanding KUs to include other forms of knowledge, such as domain-specific knowledge units. A promising approach could involve analyzing project descriptions from README files and developer discussions in pull requests or issue trackers, to extract relevant domain knowledge. By leveraging large language models (LLMs), this process can be automated and made scalable, enabling the identification and summarization of domain-specific concepts and their relationships to the software system. For example, Retrieval-Augmented Generation (RAG) architectures can be utilized, where LLMs are combined with retrieval mechanisms to dynamically pull relevant context from external knowledge sources, such as issues and pull request discussions. This allows the LLM to generate more accurate and contextually relevant summaries of domain knowledge (e.g., project-specific terminologies, workflows, or dependencies). We also encourage future works to integrate other form of technical knowledge related to libraries and APIs. This could include high-level libraries widely used in data science, such as pandas, as well as lower-level, such as `OpenSSL` for cryptographic operations, and `libcurl` for network communication. Expanding this scope could involve labeling issues with API domains [51], mapping analogous APIs between third-party libraries [32], and identifying API caveats or misuse patterns [30]. Incorporating these diverse forms of knowledge could provide deeper insights into code behavior and can significantly enhance the predictive power of models for identifying post-release defects.

## 6 Threats to Validity

**Construct validity.** We use the Java certification exams as guidelines to identify KUs because certification exams cover important topics of the Java programming language. We only focus on the core functionalities of the programming language that are listed in the topics of these exams. The set of KUs that we conceive in this paper is thus limited and does not cover all the knowledge embedded in the constructs and APIs of the Java programming language.

However, these core functionalities are covered in the Java certification exams (Java SE and Java EE) to assess the expertise of developers. We encourage future studies to refine our conceptualization of KUs and operationalize them in different ways.

Another construct validity threat stems from the fact that we did not extract KUs from the third-party libraries used by the studied systems. In this study, we only focus on those KUs that are implemented by the developers of the studied software systems. We carefully implemented the type binding resolutions using the Java JDT [11] to ensure that we are less likely to miss the KUs implemented within the studied projects.

We reused the 65 traditional metrics provided in the defect dataset, which include 54 product metrics, 5 process metrics, and 6 code ownership metrics. While there are many traditional metrics available, the selected traditional metrics for our study are widely recognized and have been extensively utilized in prior research to build and analyze defect prediction models [15, 24, 29, 42, 58].

**Internal validity.** In our study, we used a Random Forest (RF) classifier with default hyper-parameter settings as the standard configuration for building the studied models. However, model performance can vary significantly across different classifiers and hyperparameter settings since each classifier has its unique strengths and weaknesses. To address this concern, we analyzed the impact of different classifiers and hyper-parameter tuning on model performance. We experimented with a range of classifiers, including traditional methods (e.g., Decision Trees, KNN) and advanced techniques (e.g., Light-GBM, XGBoost). Our results showed that RF with default settings consistently achieved the best performance, minimizing bias associated with classifier selection and hyper-parameter tuning. Details of this analysis are provided in C.

The results reported in this paper are based on models that were trained without applying data rebalancing techniques, such as SMOTE [6]. To investigate whether class rebalancing techniques could improve model performance, we applied SMOTE and re-ran our evaluations. However, we observed no improvement in performance. Moreover, such class rebalancing techniques are not recommended when building explainable models, as they can introduce issues like concept drift and hinder model interpretability [54]. Therefore, we do not use SMOTE in our study.

**External validity.** In this study, we only considered Java projects. Future studies should broaden the scope of our study and investigate how our findings apply to software projects written in other popular programming languages. While the findings demonstrate cost-effectiveness within the context of the studied dataset, it is unclear whether this cost-effectiveness can be generalized to other defect datasets. Differences in characteristics such as dataset size, structure, or domain could influence the applicability of the results to broader contexts. However, it is worth noting that the defect dataset used in this study

consists of multiple releases of eight well-maintained and real-world software systems, ensuring a realistic evaluation.

## 7 Related Work

A number of studies have focused on analyzing different metrics for defect prediction models [1, 4, 5, 13, 20, 25, 28, 38, 40, 43, 59, 60]. In this section, we highlight a few of these studies that are more closely related to our study.

The simplest, easiest to extract, and most commonly used product metric for bug prediction is the lines of code (LOC). Zhang [59] investigated the relationship between LOC and bug in Eclipse and several NASA projects. The author observed that typical classification techniques based on LOC were able to predict defective components reasonably well. Koru et al. [26, 27] observed a power-law relationship between LOC and bug-proneness, with the latter increasing at a slower rate. They also observed that smaller modules were proportionally more bug-prone compared with larger ones. Although LOC correlates with the number of bugs [14, 47, 59], some studies find no strong evidence that LOC is a good indicator of bugs [3, 16]. Many researchers investigated other product metrics based on code, such as object-oriented metrics and complexity metrics for predicting bugs. These metrics are commonly referred to as code metrics. Gyimothy et al. [20] used eight object-oriented metrics and traditional code-size metrics (e.g., LOC) to detect bug-proneness in the Mozilla source code. Their results indicated that all code metrics (except the one that calculates the number of direct descendants of a class) were significant for bug prediction. Nagappan et al. [40] built prediction models using 18 complexity metrics along three dimensions (module, function, and class). Their evaluation results showed that complexity metrics were able to predict post-release defects with good accuracy, but no single set of metrics was applicable to all projects.

Many prior studies build bug-prediction models with process metrics and investigate the performance of the bug prediction task between models that are built with product metrics and process metrics. Nagappan et al. [39] investigated the utility of using process and product metrics to estimate post-release bug-proneness in Windows XP-SP1 and Windows 2003 Server operating systems. They observed that process and product metrics could be used to identify bug-prone modules at statistically significant levels. Moser et al. [38] proposed 18 process metrics based on the change history of the files (i.e., change metrics). They studied the performance of bug prediction models built with these 18 process metrics and 31 traditional static code metrics. They observed that predictors based on process metrics outperformed predictors based on static code metrics. In contrast with the file-level bug prediction approach, Giger et al. [19] explored bug prediction models at the method level. They developed prediction models using a combination of 15 process metrics and 9 source code metrics to assess their performance on 21 Java projects. Their findings indicated that process metrics alone were effective predictors and sig-

nificantly outperformed source code metrics. However, combining both types of metrics did not enhance the models' classification performance. Rahman and Devanbu [48] built defect prediction models using 54 code metrics and 14 process metrics and evaluated their performance on 85 releases of 12 major open-source projects. The results of their study revealed that code metrics were generally less effective than process metrics.

Majumder et al. [34] carried out an extensive study analyzing 722K commits from 700 GitHub projects to compare the effectiveness of bug prediction models that use process metrics against those that use product metrics. Their findings indicated that process metrics are superior in predicting defects. However, Dalla Palma et al. [10] observed that product metrics are better than process metrics for identifying bugs in Infrastructure-as-Code (IaC) scripts, which are used in DevOps to manage and provision infrastructure through machine-readable files.

In this paper, we empirically study knowledge units (KUs) of programming languages, a new perspective for analyzing the source code of software systems. Our goal is to understand if we can leverage KUs for predicting post-release defects. Our findings indicate that KUs are good candidates for identifying defect-prone code and offer insights into the occurrence of defects or the context in which those defects happen.

## 8 Conclusion

In this paper, we present an empirical study on how to leverage KUs of a programming language to predict post-release defects. We develop a defect prediction model called KUM, which uses KU features. We compare the performance of KUM with that of models built with traditional metrics (PROC, PROD, OWN, and TM). KUM significantly outperforms PROD, PROC and OWN. Combining KUs and traditional metrics leads to an even a higher performing model (KUM+TM) that significantly outperforms both TM and KUM. We also observe that a cost-effective model (COST_EFF) that combines the top features from KUM and TM yields a performance that is similar to KUM+TM.

Our encouraging results highlights the relevance of KUs in the context of defect prediction. Given our promising findings from this study, we encourage researchers to further explore and expand on the concept of KUs. This includes developing more advanced conceptualizations and elicitation of KUs, as well as evaluating whether the insights from our study can be generalized to other software systems and programming languages. More generally, we believe that KUs provide a brand new lens for analyzing software systems and there is still much to uncover about the potential of KUs. We thus also encourage future research to explore diverse application areas where KUs may provide valuable insights, including software quality and maintenance.

## Declarations

### Data Availability Statement (DAS)

A supplementary material package is provided online in the following link: `http://bit.ly/3GoSmHL`. The contents will be made available on a public GitHub repository once the paper is accepted.

### Funding and/or Conflicts of interests/Competing interests

The authors declared that they have no conflict of interest.

## References

1. Aggarwal KK, Singh Y, Kaur A, Malhotra R (2009) Empirical Analysis for Investigating the Effect of Object-Oriented Metrics on Fault Proneness: A Replicated Case Study. Software process: Improvement and practice 14(1):39–62
2. Ahasanuzzaman M, Oliva GA, Hassan AE (2024) Using knowledge units of programming languages to recommend reviewers for pull requests: an empirical study. Empirical Software Engineering 29(1):33
3. Andersson C, Runeson P (2007) A replicated quantitative analysis of fault distributions in complex software systems. IEEE transactions on software engineering 33(5):273–286
4. Briand LC, Wüst J, Daly JW, Porter DV (2000) Exploring the Relationship between Design Measures and Software Quality in Object-Oriented Systems. Journal of system software 51(3):245–273
5. Burrows R, Ferrari FC, Lemos OA, Garcia A, Taiani F (2010) The impact of coupling on the fault-proneness of aspect-oriented programs: An empirical study. In: Proceedings of the 21st International Symposium on Software Reliability Engineering, pp 329–338
6. Chawla NV, Bowyer KW, Hall LO, Kegelmeyer WP (2002) Smote: synthetic minority over-sampling technique. Journal of artificial intelligence research 16:321–357
7. Chen D, Chen X, Li H, Xie J, Mu Y (2019) Deepcpdp: Deep learning based cross-project defect prediction. IEEE Access 7:184832–184848
8. Chidamber SR, Kemerer CF (1994) A Metrics Suite for Object Oriented Design. IEEE Transactions on Software Engineering 20(6):476–493
9. CRAN (2020) shapper: Wrapper of Python Library 'shap' . `https://cran.r-project.org/web/packages/shapper/index.html`, (Last accessed: December 2024)
10. Dalla Palma S, Di Nucci D, Palomba F, Tamburri DA (2021) Within-project defect prediction of infrastructure-as-code using product and process metrics. IEEE Transactions on Software Engineering 48(6):2086–2104
11. Eclipse (2020) Eclipse Java development tools (JDT) . `http://www.eclipse.org/jdt/`, (Last accessed: May 2021)
12. Efron B (1983) Estimating the error rate of a prediction rule: improvement on cross-validation. Journal of the American statistical association 78(382):316–331
13. Emam KE, Melo W, Machado JC (2001) The Prediction of Faulty Classes Using Object-Oriented Design Metrics. Journal of system software 56:63–75
14. English M, Exton C, Rigon I, Cleary B (2009) Fault detection and prediction in an open-source software project. In: Proceedings of the 5th International Conference on Predictor Models in Software Engineering, pp 1–11
15. Esteves G, Figueiredo E, Veloso A, Viggiato M, Ziviani N (2020) Understanding machine learning software defect predictions. Automated Software Engineering 27(3):369–392
16. Fenton NE, Ohlsson N (2000) Quantitative analysis of faults and failures in a complex software system. IEEE Transactions on Software engineering 26(8):797–814

17. Fu M, Tantithamthavorn C (2022) Linevul: A transformer-based line-level vulnerability prediction. In: Proceedings of the 19th International Conference on Mining Software Repositories, pp 608–620
18. Ghotra B, McIntosh S, Hassan AE (2015) Revisiting the Impact of Classification Techniques on the Performance of Defect Prediction Models. In: Proceedings of the 37th IEEE International Conference on Software Engineering, pp 789–800
19. Giger E, D'Ambros M, Pinzger M, Gall HC (2012) Method-level bug prediction. In: Proceedings of the ACM-IEEE international symposium on Empirical software engineering and measurement, pp 171–180
20. Gyimothy T, Ferenc R, Siket I (2005) Empirical Validation of Object-Oriented Metrics on Open Source Software for Fault Prediction. IEEE Transactions on Software Engineering 31(10):897–910
21. Hoang T, Dam HK, Kamei Y, Lo D, Ubayashi N (2019) DeepJIT: An End-to-End Deep Learning Framework for Just-in-Time Defect Prediction. In: Proceedings of the 16th International Conference on Mining Software Repositories, pp 34–45
22. Jiarpakdee J, Tantithamthavorn C, Treude C (2018) Autospearman: Automatically mitigating correlated software metrics for interpreting defect models. In: Proceedings of the 37th International Conference on Software Maintenance and Evolution, pp 92–103
23. Jiarpakdee J, Tantithamthavorn C, Hassan AE (2019) The impact of correlated metrics on the interpretation of defect models. IEEE Transactions on Software Engineering 47(2):320–331
24. Jiarpakdee J, Tantithamthavorn C, Treude C (2020) The Impact of Automated Feature Selection Techniques on the Interpretation of Defect Models. Empirical Software Engineering 25(5):3590–3638
25. Kanmani S, Uthariaraj VR, Sankaranarayanan V, Thambidurai P (2007) Object-Oriented Software Fault Prediction Using Neural Networks. Information and software technology 49(5):483–492
26. Koru AG, Emam KE, Zhang D, Liu H, Mathew D (2008) Theory of relative defect proneness: Replicated studies on the functional form of the size-defect relationship. Empirical Software Engineering 13:473–498
27. Koru AG, Zhang D, El Emam K, Liu H (2008) An investigation into the functional form of the size-defect relationship for software modules. IEEE Transactions on Software Engineering 35(2):293–304
28. Kumar L, Misra S, Rath SK (2017) An empirical analysis of the effectiveness of software metrics and fault prediction model for identifying faulty classes. Computer Standards & Interfaces 53:1–32
29. Laaber C, Basmaci M, Salza P (2021) Predicting unstable software benchmarks using static source code features. Empirical Software Engineering 26(6):1–53
30. Li H, Li S, Sun J, Xing Z, Peng X, Liu M, Zhao X (2018) Improving api caveats accessibility by mining api caveats knowledge graph. In: 2018 IEEE International Conference on Software Maintenance and Evolution (ICSME), IEEE, pp 183–193
31. Li J, He P, Zhu J, Lyu MR (2017) Software Defect Prediction via Convolutional Neural Network. In: Proceedings of the 20th IEEE International Conference on Software Quality, Reliability and Security, pp 318–328
32. Liu Y, Liu M, Peng X, Treude C, Xing Z, Zhang X (2020) Generating concept based api element comparison using a knowledge graph. In: Proceedings of the 35th IEEE/ACM International Conference on Automated Software Engineering, pp 834–845
33. Long JD, Feng D, Cliff N (2003) Ordinal analysis of behavioral data. Handbook of psychology pp 635–661
34. Majumder S, Mody P, Menzies T (2022) Revisiting process versus product metrics: a large scale analysis. Empirical Software Engineering 27(3):60
35. Menzies T (2020) Scott Knot with nonparametric effect size and significance test . https://gist.github.com/timm/41b3a8790c1adce26d63c5874fbea393, (Last accessed: December 2024)
36. Mittas N, Angelis L (2013) Ranking and Clustering Software Cost Estimation Models through a Multiple Comparisons Algorithm. IEEE Transactions on Software Engineering 39(4):537–551
37. Molnar C (2020) Interpretable Machine Learning

38. Moser R, Pedrycz W, Succi G (2008) A Comparative Analysis of the Efficiency of Change Metrics and Static Code Attributes for Defect Prediction. In: Proceedings of the 30th International Conference on Software Engineering, p 181–190
39. Nagappan N, Ball T, Murphy B (2006) Using Historical In-Process and Product Metrics for Early Estimation of Software Failures. In: Proceedings of the 17th International Symposium on Software Reliability Engineering, pp 62–74
40. Nagappan N, Ball T, Zeller A (2006) Mining Metrics to Predict Component Failures. In: Proceedings of the 28th International Conference on Software Engineering, pp 452–461
41. Nam D, Macvean A, Hellendoorn V, Vasilescu B, Myers B (2024) Using an llm to help with code understanding. In: Proceedings of the 46th International Conference on Software Engineering, pp 1–13
42. Nuñez-Varela AS, Pérez-Gonzalez HG, Martínez-Perez FE, Souberville-Montalvo C (2017) Source code metrics: A systematic mapping study. Journal of Systems and Software 128:164–197
43. Olague HM, Etzkorn LH, Gholston S, Quattlebaum S (2007) Empirical Validation of Three Software Metrics Suites to Predict Fault-Proneness of Object-Oriented Classes Developed Using Highly Iterative or Agile Software Development Processes. IEEE Transactions on Software Engineering 33(6):402–419
44. Oracle (2022) Oracle Certified Associate, Java SE 8 Programmer. `https://education.oracle.com/oracle-certified-associate-java-se-8-programmer/trackp_333`, (Last accessed: April 2023)
45. Oracle (2022) Oracle Certified Professional, Java EE 7 Application Developer. `https://education.oracle.com/oracle-certified-professional-java-ee-7-application-developer/pexam_1Z0-900`, (Last accessed: April 2023)
46. Oracle (2022) Oracle Certified Professional, Java SE 8 Programmer. `https://education.oracle.com/oracle-certified-professional-java-se-8-programmer/trackp_357`, (Last accessed: April 2023)
47. Ostrand TJ, Weyuker EJ, Bell RM (2005) Predicting the location and number of faults in large software systems. IEEE Transactions on Software Engineering 31(4):340–355
48. Rahman F, Devanbu P (2013) How, and why, process metrics are better. In: Proceedings of the 35th International Conference on Software Engineering, pp 432–441
49. Rajbahadur GK, Wang S, Oliva GA, Kamei Y, Hassan AE (2022) The impact of feature importance methods on the interpretation of defect classifiers. IEEE Transactions on Software Engineering 48(7):2245–2261
50. Romano J, Kromrey J, Coraggio J, Skowronek J (2006) Appropriate statistics for ordinal level data: Should we really be using t-test and Cohen's d for evaluating group differences on the NSSE and other surveys? In: Annual meeting of the Florida Association of Institutional Research, pp 1–3
51. Santos F, Vargovich J, Trinkenreich B, Santos I, Penney J, Britto R, Pimentel JF, Wiese I, Steinmacher I, Sarma A, et al. (2023) Tag that issue: applying API-domain labels in issue tracking systems. Empirical Software Engineering 28(5):116
52. SciTools (2000) Understand by SciTools. `https://scitools.com`, (Last accessed: December 2024)
53. Shihab E, Jiang ZM, Ibrahim WM, Adams B, Hassan AE (2010) Understanding the Impact of Code and Process Metrics on Post-Release Defects: A Case Study on the Eclipse Project. In: Proceedings of the 4th International Symposium on Empirical Software Engineering and Measurement, pp 1–10
54. Tantithamthavorn C, Hassan AE (2018) An Experience Report on Defect Modelling in Practice: Pitfalls and Challenges. In: Proceedings of the 40th International Conference on Software Engineering: Software Engineering in Practice, p 286–295
55. Tantithamthavorn C, McIntosh S, Hassan AE, Matsumoto K (2017) An Empirical Comparison of Model Validation Techniques for Defect Prediction Models. IEEE Transactions on Software Engineering 43(1):1–18
56. Yatish S, Jiarpakdee J, Thongtanunam P, Tantithamthavorn C (2019) Mining Software Defects: Should We Consider Affected Releases? In: Proceedings of the 41st International Conference on Software Engineering, p 654–665
57. Zain ZM, Sakri S, Ismail NHA (2023) Application of deep learning in software defect prediction: systematic literature review and meta-analysis. Information and Software

Technology 158

58. Zhang F, Zheng Q, Zou Y, Hassan AE (2016) Cross-Project Defect Prediction Using a Connectivity-Based Unsupervised Classifier. In: Proceedings of the 38th International Conference on Software Engineering, p 309–320

59. Zhang H (2009) An investigation of the relationships between lines of code and defects. In: Proceedings of the 25th International Conference on Software Maintenance, pp 274–283

60. Zhou Y, Leung H (2006) Empirical Analysis of Object-Oriented Design Metrics for Predicting High and Low Severity Faults. IEEE Transactions on Software Engineering 32(10):771–789

# APPENDICES

## A Java Certification Exams and Knowledge Units

Table 6 summarizes our KUs and their capabilities.

## B Traditional code metrics

The list of process metrics employed in this study is presented in Table 7. Table 8 presents the list of process metrics and Table 9 presents the list of ownership metrics that we study in this paper. These tables were adapted from the original paper of Yatish et al. [56].

**Table 6** Knowledge Units derived from the Java SE 8 Programmer I Exam, Java SE 8 Programmer II Exam, and Java EE Developer Exam.

| Certification Exam | Knowledge Unit (KU) | Key Capabilities |
|---|---|---|
| Java SE I | **[K1]** Data Type | **[C1]** Declare and initialize different types of variables (e.g., primitive types, parameterized types, and arrays), including casting between primitive types. |
| | **[K2]** Operator & Decision | **[C1]** Use Java operators (e.g., assignment and postfix operators); use parentheses to override operator precedence.<br>**[C2]** Test equality between strings and other objects using `==` and `equals()`.<br>**[C3]** Create and use `if`, `if-else`, and ternary constructs.<br>**[C4]** Use a `switch` statement. |
| | **[K3]** Array | **[C1]** Declare, instantiate, initialize and use a one-dimensional array<br>**[C2]** Declare, instantiate, initialize and use a multi-dimensional array |
| | **[K4]** Loop | **[C1]** Create and use `while` loops<br>**[C2]** Create and use `for` loops, including the `enhanced for` loop<br>**[C3]** Create and use `do-while` loops<br>**[C4]** Use `break` statement<br>**[C5]** Use `continue` statement |
| | **[K5]** Method & Encapsulation | **[C1]** Create methods with arguments and return values<br>**[C2]** Apply the "static" keyword to methods, fields, and blocks<br>**[C3]** Create an overloaded method and overloaded constructor<br>**[C4]** Create a constructor chaining (use "this()" method to call one constructor from another constructor<br>**[C5]** Use variable length arguments in the methods<br>**[C6]** Use different access modifiers (e.g., private and protected) other than "default"<br>**[C7]** Apply encapsulation: identify set and get method to initialize any private class variables<br>**[C8]** Apply encapsulation: Immutable class generation-final class and initialize private variables through the constructor |
| | **[K6]** Inheritance | **[C1]** Use basic polymorphism (e.g., a superclass refers to a subclass)<br>**[C2]** Use polymorphic parameter (e.g., pass instances of a subclass or interface to a method)<br>**[C3]** Create overridden methods<br>**[C4]** Create "abstract" classes and "abstract" methods<br>**[C5]** Create "interface" and implement the interface<br>**[C6]** Use "super()" and the "super" keyword to access the members(e.g., fields and methods) of a parent class<br>**[C7]** Use casting in referring a subclass object to a superclass object |

| Certification Exam | Knowledge Unit (KU) | Key Capabilities |
| --- | --- | --- |
| Java SE II | **[K7]** Advanced Class Design | **[C1]** Create inner classes, including static inner classes, local classes, nested classes, and anonymous inner classes<br>**[C2]** Develop code that uses the final<br>**[C3]** Use enumerated types including methods and constructors in an "enum" type<br>**[C4]** Create singleton classes and immutable classes |
| | **[K8]** Generics & Collection | **[C1]** Create and use a generic class<br>**[C2]** Create and use `ArrayList`, `TreeSet`, `TreeMap`, and `ArrayDeque` **[C3]** Use `java.util.Comparator` and `java.lang.Comparable` interfaces<br>**[C4]** Iterate using forEach methods of List |
| | **[K9]** Functional Interface | **[C1]** Use the built-in interfaces included in the `java.util.function` packages such as `Predicate`, `Consumer`, `Function`, and `Supplier`<br>**[C2]** Develop code that uses primitive versions of functional interfaces<br>**[C3]** Develop code that uses binary versions of functional interfaces<br>**[C4]** Develop code that uses the `UnaryOperator` interface |
| | **[K10]** Stream API | **[C1]** Develop code to extract data from an object using `peek()` and `map()` methods, including primitive versions of the `map()` method<br>**[C2]** Search for data by using search methods of the Stream classes, including `findFirst`, `findAny`, `anyMatch`, `allMatch`, `noneMatch`<br>**[C3]** Develop code that uses the Optional class<br>**[C4]** Develop code that uses Stream data methods and calculation methods<br>**[C5]** Sort a collection using Stream API<br>**[C6]** Save results to a collection using the collect method<br>**[C7]** Use `flatMap()` methods in the Stream API |
| | **[K11]** Exception | **[C1]** Create a try-catch block<br>**[C2]** Use catch, multi-catch, and finally clauses<br>**[C3]** Use autoclose resources with a try-with-resources statement<br>**[C4]** Create custom exceptions and autocloseable resources<br>**[C5]** Create and invoke a method that throws an exception<br>**[C6]** Use common exception classes and categories(such as `NullPointerException`, `ArithmeticException`, `ArrayIndexOutOfBoundsException`, `ClassCastException`)<br>**[C5]** Use assertions |
| | **[K12]** Date Time API | **[C1]** Create and manage date-based and time-based events including a combination of date and time into a single object using `LocalDate`, `LocalTime`, `LocalDateTime`, `Instant`, `Period`, and `Duration`<br>**[C2]** Formatting date and times values for using different timezones<br>**[C3]** Create and manage date-based and time-based events using Instant, Period, Duration, and Temporal Unit<br>**[C4]** Create and manipulate calendar data using classes from `java.time.LocalDateTime`, `java.time.LocalDate`, `java.time.LocalTime`, `java.time.format.DateTimeFormatter`, and `java.time.Period` |
| | **[K13]** IO | **[C1]** Read and write data using the console<br>**[C2]** Use `BufferedReader`, `BufferedWriter`, `File`, `FileReader`, `FileWriter`, `FileInputStream`, `FileOutputStream`, `ObjectOutputStream`, `ObjectInputStream`, and `PrintWriter` in the `java.io` package |
| | **[K14]** NIO | **[C1]** Use the Path interface to operate on file and directory paths<br>**[C2]** Use the Files class to check, read, delete, copy, move, and manage metadata a file or directory |
| | **[K15]** String Processing | **[C1]** Search, parse and build strings<br>**[C2]** Manipulate data using the `StringBuilder` class and its methods<br>**[C3]** Use regular expression using the `Pattern` and `Matcher` class<br>**[C4]** Use string formatting |
| | **[K16]** Concurrency | **[C1]** Create worker threads using `Runnable`, `Callable` and use an `ExecutorService` to concurrently execute tasks<br>**[C2]** Use `synchronized` keyword and `java.util.concurrent.atomic` package to control the order of thread execution<br>**[C3]** Use `java.util.concurrent` collections and classes including `CyclicBarrier` and `CopyOnWriteArrayList`<br>**[C4]** Use parallel Fork/Join Framework |
| | **[K17]** Database | **[C1]** Describe the interfaces that make up the core of the JDBC API, including the `Driver`, `Connection`, `Statement`, and `ResultSet` interfaces<br>**[C2]** Submit queries and read results from the database, including creating statements, returning result sets, iterating through the results, and properly closing result sets, statements, and connections |
| | **[K18]** Localization | **[C1]** Read and set the locale by using the Locale object<br>**[C2]** Build a resource bundle for each locale and load a resource bundle in an application |

| Certification Exam | Knowledge Unit (KU) | Key Capabilities |
|---|---|---|
| Java EE | **[K19]** Java Persistence | **[C1]** Create JPA Entity and Object-Relational Mappings (ORM) |
| | | **[C2]** Use Entity Manager to perform database operations, transactions, and locking with JPA entities |
| | | **[C3]** Create and execute JPQL statements |
| | **[K20]** Enterprise Java Bean | **[C1]** Create session EJB components containing synchronous and asynchronous business methods, manage the life cycle container callbacks, and use interceptors. |
| | | **[C2]** Create EJB timers |
| | **[K21]** Java Message Service API | **[C1]** Implement Java EE message producers and consumers, including Message-Driven beans |
| | | **[C2]** Use transactions with JMS API |
| | **[K22]** SOAP Web Service | **[C1]** Create SOAP Web Services and Clients using JAX-WS API |
| | | **[C2]** Create marshall and unmarshall Java Objects by using JAXB API |
| | **[K23]** Servlet | **[C1]** Create Java Servlet and use HTTP methods |
| | | **[C2]** Handle HTTP headers, parameters, cookies |
| | | **[C3]** Manage servlet life cycle with container callback methods and WebFilters |
| | **[K24]** Java REST API | **[C1]** Apply REST service conventions |
| | | **[C2]** Create REST Services and clients using JAX-RS API |
| | **[K25]** Websocket | **[C1]** Create WebSocket Server and Client Endpoint Handlers |
| | | **[C2]** Produce and consume, encode and decode WebSocket messages |
| | **[K26]** Java Server Faces | **[C1]** Use JSF syntax and use JSF Tag Libraries |
| | | **[C2]** Handle localization and produce messages |
| | | **[C3]** Use Expression Language (EL) and interact with CDI beans |
| | **[K27]** Contexts and Dependency Injection (CDI) | **[C1]** Create CDI Bean Qualifiers, Producers, Disposers, Interceptors, Events, and Stereotypes |
| | **[K28]** Batch Processing | **[C1]** Implement batch jobs using JSR API |

**Table 7** The list of traditional code metrics that we select for this study.

| Granularity | Metric Name | Understand (API) name | Definition |
|---|---|---|---|
| File LeveL(37) | Average Cyclomatic Complexity | AvgCyclomatic | Average cyclomatic complexity for all nested functions or methods. |
| | Average Modified Cyclomatic Complexity | AvgCyclomaticModified | Average modified cyclomatic complexity for all nested functions or methods. |
| | Average Strict Cyclomatic Complexity | AvgCyclomaticStrict | Average strict cyclomatic complexity for all nested functions or methods. |
| | Average Essential Cyclomatic Complexity | AvgEssential | Average Essential complexity for all nested functions or methods. |
| | Average Number of Lines | AvgLine | Average number of lines for all nested functions or methods. |
| | Average Number of Blank Lines | AvgLineBlank | Average number of blanks for all nested functions or methods. |
| | Average Number of Lines of Code | AvgLineCode | Average number of lines containing source code for all nested functions or methods. |
| | Average Number of Lines with Comments | AvgLineComment | Average number of lines containing comments for all nested functions or methods. |
| | Classes | CountDeclClass | Number of classes. |
| | Class Methods | CountDeclClassMethod | Number of class methods. |
| | Class Variables | CountDeclClassVariable | Number of class variables. |
| | Function | CountDeclFunction | Number of functions. |
| | Instance MethodS(NIM) | CountDeclInstanceMethod | Number of instance methods. |
| | Instance VariableS(NIV) | CountDeclInstanceVariable | Number of instance variables. |
| | Local Methods | CountDeclMethod | Number of local methods. |
| | Local Default Visibility Methods | CountDeclMethodDefault | Number of local default methods. |
| | Private MethodS(NPM) | CountDeclMethodPrivate | Number of local private methods. |
| | Protected Methods | CountDeclMethodProtected | Number of local protected methods. |
| | Public Methods | CountDeclMethodPublic | Number of local public methods. |
| | Physical LineS(NL) | CountLine | Number of all lines. |
| | Blank Lines of Code (BLOC) | CountLineBlank | Number of blank lines. |
| | Source Lines of Code | CountLineCode | Number of lines containing source code. [aka LOC] |
| | Declarative Lines of Code | CountLineCodeDecl | Number of lines containing declarative source code. |
| | Executable Lines of Code | CountLineCodeExe | Number of lines containing executable source code. |
| | Lines with Comments | CountLineComment | Number of lines containing comments. |
| | Semicolons | CountSemicolon | Number of semicolons. |
| | Statements | CountStmt | Number of statements. |
| | Declarative Statements | CountStmtDecl | Number of declarative statements. |
| | Executable Statements | CountStmtExe | Number of executable statements. |
| | Max Cyclomatic Complexity | MaxCyclomatic | Maximum cyclomatic complexity of all nested functions or methods. |
| | Max Modified Cyclomatic Complexity | MaxCyclomaticModified | Maximum modified cyclomatic complexity of nested functions or methods. |
| | Max Strict Cyclomatic Complexity | MaxCyclomaticStrict | Maximum strict cyclomatic complexity of nested functions or methods. |
| | Comment to Code Ratio | RatioCommentToCode | Ratio of comment lines to code lines. |
| | Sum Cyclomatic Complexity (WMC) | SumCyclomatic | Sum of cyclomatic complexity of all nested functions or methods. |
| | Sum Modified Cyclomatic Complexity | SumCyclomaticModified | Sum of modified cyclomatic complexity of all nested functions or methods. |
| | Sum Strict Cyclomatic Complexity | SumCyclomaticStrict | Sum of strict cyclomatic complexity of all nested functions or methods. |
| | Sum Essential Complexity | SumEssential | Sum of essential complexity of all nested functions or methods |
| Class LeveL(5) | Coupling Between ObjectS(CBO) | CountClassCoupled | Number of other classes coupled to. |
| | Number of ChildreN(NOC) | CountClassDerived | Number of immediate subclasses. |
| | Depth of Inheritance Tree (DIT) | MaxInheritanceTree | Maximum depth of class in inheritance tree. |
| | Lack of Cohesion in MethodS(LCOM) | PercentLackOfCohesion | Lack of Cohesion in Methods |
| | Response for clasS(RFC) | CountDeclMethodAll | Number of methods, including inherited ones. |
| Method LeveL(12) | InputS(FANINN) | CountInput {Min, Median, Max} | Number of calling subprograms plus global variables read. |
| | OutputS(FANOUT) | CountOutput {Min, Median, Max} | Number of called subprograms plus global variables set. |
| | PathS(NPATH) | CountPath {Min, Median, Max} | Number of possible paths, not counting abnormal exits or gotos. |
| | Nesting | MaxNesting {Min, Median, Max} | Maximum nesting level of control constructs. |

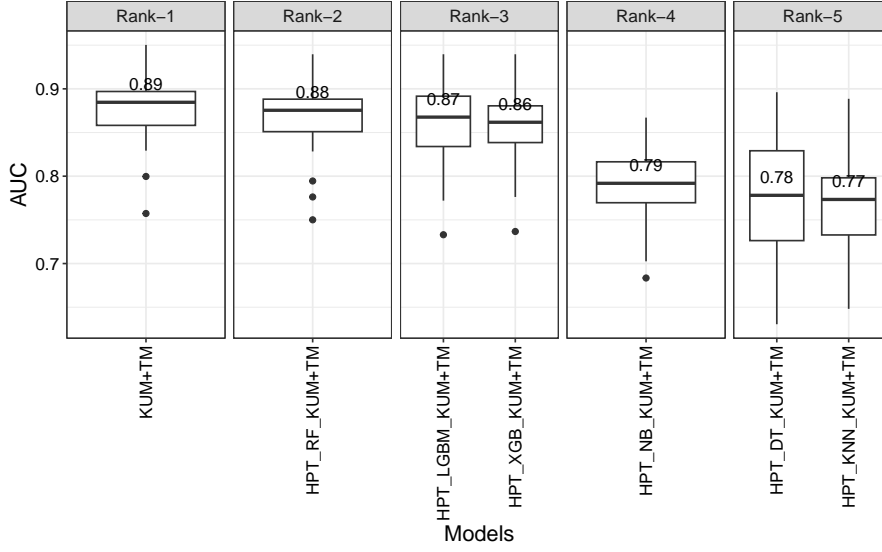**Table 8** The list of process metrics that we select for this study.

| Metric Name | Description |
| --- | --- |
| COMM | The number of Git commits |
| ADDED_LINES | The normalized number of lines added to the file |
| DEL_LINES | The normalized number of lines deleted from the file |
| ADEV | The number of active developers |
| DDEV | The number of distinct developers |

**Table 9** The list of code ownership metrics that we select for this study.

| Metric Name | Description |
| --- | --- |
| MINOR_COMMIT | The number of unique developers who have contributed less than 5% of the total code changes (i.e., Git commits) on the file |
| MINOR_LINE | The number of unique developers who have contributed less than 5% of the total lines of code on the file |
| MAJOR_COMMIT | The number of unique developers who have contributed more than 5% of the total code changes (i.e., Git commits) on the file |
| MAJOR_LINE | The number of unique developers who have contributed more than 5% of the total lines of code on the file |
| OWN_COMMIT | The proportion of code changes (i.e., Git commits) made by the developer who has the highest contribution of code changes on the file |
| OWN_LINE | The proportion of lines of code written by the developer who has the highest contribution of lines of code on the file developers |

**Table 10** The studied classification algorithms with different parameter settings for hyper parameter optimization.

| Classifier Name | Classification Algorithm | Parameter Name | Parameter Description | Studied candidate parameter values |
| --- | --- | --- | --- | --- |
| KNN | K-Nearest Neighbour | n_neighbors | The number of neighbors required for each sample | {1, 5, 9, 13, 17, 20} |
| NB | Naive Bayes | var_smoothing | The portion of the largest variance of all features that is added to variances for calculation stability. | {1e-5, 1e-9, 1e-11, 1e-15} |
| DT | Decision Tree | criterion | The function to measure the quality of a split. | {'gini', 'entropy', 'log_loss'} |
| | | max_depth | The maximum depth of the tree. | {None, 5, 10} |
| | | ccp_alpha | Complexity parameter used for Minimal Cost-Complexity Pruning. | {0.0001, 0.001, 0.01, 0.1, 0.5} |
| RF | Random Forest | n_estimator | The number of trees in the forest. | {10, 50, 100, 200} |
| | | max_depth | The maximum depth of the tree. If None, then nodes are expanded until all leaves are pure. | {None, 5, 10} |
| XGB | XGBoost | n_estimator | The number of boosting rounds or trees to build. | {10, 50, 100, 200} |
| | | max_depth | The maximum depth of the tree. | {None, 5, 10} |
| | | learning_rate | This parameter controls how much the model is adjusted in response to the estimated error each time the model weights are updated. | {0.1, 0.01, 0.001} |
| LGBM | LightGBM | n_estimator | The number of boosted trees to fit. | {10, 50, 100, 200} |
| | | num_leaves | The maximum tree leaves for base learners. | {None, 5, 10} |
| | | learning_rate | The Boosting learning rate that controls how quickly the model adjusts to the error in each iteration of training | {0.1, 0.01, 0.001} |

**Fig. 11** The AUC distribution of hyper-parameter-tuned models and our original KUM+TM (random forest with default hyper-parameter values). The names of hyper-parameter tuning models begin with "HPT_X", where X represents the name of one of the classifier we studied. The models are grouped based on their performance rankings determined by the Scott-Knott ESD (SK-ESD) method, where a lower SK-ESD rank indicates a better-performing model.

## C On the influence of classifier choice and hyper-parameter tuning

To evaluate the impact of different classifiers and hyper-parameter tuning on the performance of KUM+TM, we experiment with a diverse set of classifiers: Naive Bayes (NB), K-Nearest Neighbor (KNN), Decision Trees (DT), Random Forest (RF), XGBoost (XGB), and LightGBM (LGBM). XGB and LGBM represent state-of-the-art classifiers known for their advanced capabilities, while the other classifiers are well-established and widely used in empirical software engineering research.

   We employ the out-of-sample bootstrap model validation technique with 100 repetitions to ensure robust evaluation. In each iteration, a bootstrap sample is generated to train the model, and the model is tested on the out-of-sample data (i.e., data not included in the bootstrap sample). To identify the optimal configuration for each classifier, we use Scikit-learn's GridSearchCV[7]. This method performs an exhaustive search over a predefined grid of hyper-parameters (detailed in Table 10) and employs 10-fold cross-validation to evaluate each configuration. GridSearchCV is applied to the bootstrap sample data, identifying the model configuration that achieves the best performance

---

[7] `https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html`

based on hyper-parameter tuning. The best-tuned model is then tested on the out-of-sample data.

For model evaluation, we follow the methodology described in Step 1.4 of Section 4.1, ensuring consistency in assessing the effectiveness of the classifiers and their tuned configurations.

Figure 11 presents the AUC distribution of KUM+TM and the hyper-parameter-tuned variations. The KUM+TM with the default parameter settings is the top ranked one compared to the other hyper-parameter-tuned variations. The hyper-parameter-tuned random forest model (HPT_RF_KUM+TM) ranks the second achieving a median AUC of 0.88 which is lower than KUM+TM (the median AUC of KUM+TM is 0.89). In contrast, models built using traditional classifiers such as Naive Bayes (NB), Decision Trees (DT), and K-Nearest Neighbor (KNN) exhibit suboptimal performance, with median AUC scores below 0.80. Of these, the KNN-based model performs the worst, with a median AUC below 0.76. Advanced classifiers like XGBoost and LightGBM achieve better results, with median AUC scores exceeding 0.85, yet they still fail to surpass the performance of the original KUM+TM model. Thus, simply switching to a different classifier with tuned hyper-parameters does not improve the performance of KUM+TM.

**Summary**

The performance of our original KUM+TM built with the default parameter settings with RF classifier ranks the top performing model among the models built with different classifiers and tuned hyper-parameters. Thus, our original KUM+TM does not improve by simply using a different classification algorithm with tuned hyper-parameters.