

RESEARCH ARTICLE

Security-based code smell definition, detection, and impact quantification in Android

Yi Zhong | Mengyu Shi | Jiawei He | Chunrong Fang | Zhenyu Chen*

¹ State Key Laboratory for Novel Software Technology, Nanjing University 210093, China, Nanjing

Correspondence

Zhenyu Chen, Software Institute, Nanjing University 210093, Nanjing, China.
Email: zychen@nju.edu.cn

Funding information

National natural science foundation of China (Key Program): Software testing technology for security-critical deep learning system: 61832009; National natural science foundation of China (Key Program): Data-driven testing methodologies and echnologies for intelligent software systems: 61932012

Abstract

Android occupies a high market share, and its broad functions make Android security matter. Research reveals that many security issues are caused by insecure coding practices. As a poor design indicator, code smell threatens the safety and quality assurance of Android applications (apps). Although previous works revealed specific problems associated with code smells, the field still lacks research reflecting Android features. Moreover, the cost and time limit developers to repairing numerous smells timely. We conducted a study, including definition, detection, and impact quantification for Android code smell (DefDIQ): (1) define 15 novel code smells in Android from a security programming perspective; meanwhile, we provide suggestions on how to eliminate or mitigate them; (2) implement DACS to automatically detect the custom code smells based on ASTs; (3) investigate the correlation between individual smells with DACS detection results, and select suitable code smells to construct fault counting models, then quantify their impact on quality, and thereby generating code smell repair priorities. We conducted experiments on 4,575 open-source apps, and the findings are: (i) Lin's CCC between DACS and manual detection results reaches 0.9994, verifying the validity; (ii) the fault counting model constructed by ZINB is superior to NB (AIC = 517.32, BIC = 522.12); some smells do indicate fault-proneness, and we identify such avoidable poor designs; (iii) different code smells have different importance and the repair priorities constructed provide a practical guideline for researchers and inexperienced developers.

KEYWORDS:

Android code smell, quality assurance, security, repair priority

1 | INTRODUCTION

The widespread popularity of Android smartphones and other mobile devices makes the app market dynamic and challenging. Go-anywhere applications support a wide range of financial, social, and enterprise services for any user with a cellular data plan. During the ongoing improvement of Android, various software patterns and poor coding practices received increasing security and quality concerns. Once apps' security issues arise, the privacy and security of uncountable users will be at risk. The security and quality assurance of Android apps become crucial and vital to keep appealing and adapting to new devices. As a metric indicating the sub-optimal design, smells are the main culprit^{1,2}.

⁰ **Abbreviations:** AST, Abstract Syntax Tree; CCC, Concordance Correlation Coefficient; NB, Negative Binomial; ZINB, Zero-Inflated Negative Binomial

By statistics, Android occupies 71.55% market share, about 2.5 times more than IOS³. The high market share makes Android more vulnerable to attacks by external sources, and security risks from poor programming selections within the code will expose users, and the platform, to varying levels of risk⁴. In this context, code smell emerged as a research hit to identify potentially poor design⁵. Although many works revealed the specific issues related to code smell, research on the definition and identification of Android remains scarce⁶. Compared with traditional software, code smells are distributed differently in Android, and the smells that truly reflect Android characteristics have not been widely studied. Most researchers following the definitions of Fowler et al.⁷, and Brown et al.⁸, there is an absence of new literature that identifies and defines Android code smells. Consequently, it becomes challenging and necessary to define and detect more code smells related to security and quality assurance.

In practical development, massive code smells occur attributed to lousy coding habits and poor design patterns. Some developers prefer to refactor the code when performance issues arise rather than anticipate and repair smells in advance⁹. The costs and time hinder developers from repairing all code smells before they threaten the app's security, quality, and hardware performance. Research finds that some code smells do indicate fault-prone code, different smells may potentially affect the apps differently, and arbitrary refactoring some smells may increase fault-proneness in some cases¹⁰. Unfortunately, refactoring code smells without guidance may be a waste of time as repairing the ones with lower importance first. The maintenance cost of apps will increase exponentially with the detection and repair time of smells. Especially after release, Android's increasing complexity and scale will substantially increase the labor and costs.

Motivated by these observations, we researched the definition, detection, and static analysis of the importance level for Androids code smell. Specifically, we defined 15 novel Android code smells, marked the parts of the code vulnerable to security risks, and provided suggestions for elimination or mitigation. No existing tool can identify the 15 smells, so we developed our own detection tool—DACS (Detect Android Code Smell) based on AST to identify whether the code smells occur in a Method body. Besides, we established a dataset of 4,575 open-source apps and collected fault information from Github by version control and keywords match. According to what we informed, there is no work to (i) a large-scale survey to collect faults in Android apps, (ii) analyze the impact of Android code smells on the fault occurrence, and (iii) quantify the code smells impact degree and generate repair priority. To obtain the code smell importance level, we investigated the distribution of code smells and fault data in the apps. We constructed fault counting models based on NB and ZINB, where the arguments are the occurrence of code smells, and the dependent variable is the fault counts in apps. Static analysis of the effect of code smells on the magnitude of faults by regression models. Then we filtered smells with low-quality impact by a two-sided hypothesis test to determine the final arguments. By the regression coefficient of the model, we generate the repair priority for the code smells.

DefDIQ is dedicated to reducing security risks, lowering technical debt, and improving code cleanliness. The fault-centric code smell repair priority provides insight into the quality and security assurance of Android apps. This work helps guide developers to identify common security structure issues and promulgate the impact of programming choices in creating secure apps. The main contributions are as follows.

- **Definition.** We defined 15 novel Android code smells from the perspective of coding specifications, programming practices, and code features related to Android security and found the relationship between smells and faults, which will assist developers in discovering more poor designs and obtaining high-quality, low-risk apps.
- **Tool.** We developed DACS to automatically identify targeted code representing the custom-defined smells by establishing smells rule library. And Lin's CCC between DACS and manual detection results reached 0.994, verifying its effectiveness.
- **Repair priority.** We constructed fault counting models to quantify the importance of code smells by taking faults as a carrier with NB and ZINB. First, we find ZINB is a better modeling technique for investigating the relationship between faults data and code smells (compared to NB). Second, the findings reveal an interesting relationship between code smells and faults: distinct smells affect faults differently, some are in combination effects and some barely affect at all. Third, we quantify 15 code smells and generate repair priority, which provides a practical guideline for inexperienced developers to remediate smells timely and enhance the benefits.

The remainder is organized as follows. Section 2 presents the background and related knowledge. Section 3 explains DefDIQ method. Section 4 describes the experimental setup, including the data set and evaluation metrics. Section 5 discusses the evaluation and the results of the conducted experiment. Section 6 presents threats to validity, while section 7 introduces the related works. The final section concludes the paper and provides directions for future work.

2 | BACKGROUND

2.1 | Code Smell in Android

Code smell was originally proposed by Kent Bech on Wiki and popularized by Fowler⁷. Sobrinho et al.¹¹ investigated plentiful works related to code smell from 1990 to 2017 and found the ones defined by Fowler were extensively studied. Code smell is a poor, sub-optimal coding structure, which may hinder comprehension¹², increase security risks and fault-proneness, and reduce maintainability¹³. When extending the application with new functions, the subtle latent errors may bring fatal consequences and hinder the project's progress. Some fields, such as aerospace, medical, self-driving, etc., will suffer huge losses and irreversible results upon software faults¹⁴.

Most Android apps are developed in Java/C++, so previous studies identified negative designs based on Java static code metrics^{15,16}. In essence, these studies still work for Android apps, but Android features distribute smells differently¹⁷, and some specific code smells are more frequent¹⁸. The distinguishing features of Android app development, as compared to traditional software, are the highly integrated environment and heavy reliance on package usage. Therefore the constant permission acquisition and confusing API calls during usage pose a security risk. The studies revealed that security risks are often accompanied by irregular coding structures, and Reimann et al.¹⁹ summarise a set of poor programming habits (e.g., a non-static internal class containing a reference to an external class), namely Android-specific code smells. The Android code smells may threaten the security, data integrity, and software quality of mobile applications^{2,20}. Given the above research, this paper defines security-related code smells from the perspective of Android security, i.e., a symptom of code predicting security risks.

2.2 | Counting Model

In this paper, we utilize faults to comprehensively assess the Android quality, including software defects such as performance, memory, and security. The importance of smells is quantified by examining the distribution relationship between faults and code smells, where faults and smells are obtained by counting. The counting model is the simplest way to fit such a series of integers, where the count type is the response variable. The Poisson regression model is the classical count model, but with the precondition: the mean = the variance. In addition, the Poisson distribution requires the observed values to be independent. If a data set is interdependent, the correlation between the data must be considered. Gupta et al.²¹ considered that the Poisson distribution needs to satisfy its conditions; otherwise, the parameter estimates will heavily deviate. To permit the variance of the observed data over the mean, Greenwood et al.²² extended the Poisson regression to a NB regression. The NB distribution consists of a connected compound Poisson distribution, with the Poisson mean obeying the γ distribution, and its probability distribution equation is:

$$P_r(Y = y) = \frac{\Gamma(y + \tau)}{y! \Gamma(\tau)} \left(\frac{\tau}{\lambda + \tau} \right)^\tau \left(\frac{\lambda}{\lambda + \tau} \right)^y, \quad y = 0, 1, \dots; \lambda, \tau > 0 \quad (1)$$

where $\lambda = E(Y)$, and τ refers to the over-dispersion. Y denotes the counting type argument with non-negative integers, and the variance is $\lambda + \lambda^2/\tau$. When $\tau \rightarrow \infty$, the variance equals the mean, and the data degenerates to Poisson distribution.

In some cases, the probability of a "zero" event is easily underestimated, making it impossible to fit the two data distributions satisfactorily. Johnson et al.²³ observed the "zero inflation" phenomenon and proposed the fence model to overcome this shortcoming. The fence model consists of Poisson-Hurdle model (PH) and Negative Binomial Hurdle (NBH), namely, assuming the zero and positive variables in the response variable originate from separate data distributions. Diane Lambert et al.²⁴ proposed a mixed distribution model called zero-inflated Poisson (ZIP). Then, Greene et al.²⁵ applied the principle to the NB distribution and proposed ZINB regression, where the Berndt-Hall-Hausman (BHHH) method is utilized to calculate the standard errors of the model parameters²⁶. The ZINB regression model separates the data into two components. The variables above zero obey the NB distribution, while the zero variables obey the discrete distribution, with the following equation:

$$P_r(Y = y) = \begin{cases} p + (1 - p)(1 + \frac{\lambda}{\tau})^{-\tau}, & y = 0 \\ (1 - p) \frac{\Gamma(y + \tau)}{y! \Gamma(\tau)} (1 + \frac{\lambda}{\tau})^{-\tau} (1 + \frac{\tau}{\lambda})^{-y}, & y = 1, 2, \dots \end{cases} \quad (2)$$

where, τ represents the overdispersion, and the mean and variance are: $E(Y) = (1 - p)\lambda$ and $var(Y) = (1 - p)\lambda(1 + p\lambda + \lambda/\tau)$. When $\tau \rightarrow \infty$ and $p \rightarrow 0$, the data obey ZINB and NB distributions, respectively.

To determine the suitable analysis technology, we test the distribution of our data. Since the faults and smells data studied in this paper do not obey normal distributions and nearly half of the data were zero, see Fig. 3, we adopted the NB and ZINB models to investigate the distribution relationship between the data.

2.3 | Spearman Rank Correlation Coefficient²⁷

There are two classical coefficients to measure the correlation between indicators: Spearman rank correlation coefficient and Pearson correlation coefficient²⁸. The Pearson correlation coefficient has two restrictions: (1) the data obey the normal distribution; (2) the data units are consistent, and the zeros are relative, not absolute. If the measured metrics do not meet the Pearson conditions, it is necessary to consider the Spearman rank correlation coefficient. It tends to be applied for calculating the correlation between the sorted variables, so the initial data need to be transformed into ordered data. Assuming x_i and y_i are a pair of order data with a monotonically increasing trend, the Spearman rank correlation coefficient is calculated as follows:

$$\rho = \frac{\sum i(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum i(x_i - \bar{x})^2 \sum i(y_i - \bar{y})^2}} \quad (3)$$

Table 1 gives the correlation coefficients and the correlation levels. When the coefficient is 0, it represents that the two metrics are entirely unrelated; when the coefficient is 1, it denotes that the two are entirely related.

Table 1 Spearman Rank Correlation Coefficient

Correlation coefficient	Correlation level
0.0-0.1	No correlation
0.1-0.3	Weak correlation
0.3-0.5	Moderate correlation
0.5-0.7	Medium-high correlation
0.7-0.9	High correlation
0.9-1.0	Full correlation

The Spearman correlation is not strict on data preconditions and ignores the overall data distribution and sample size. This paper uses Spearman's rank correlation coefficient to measure the correlation between code smells.

3 | DEFDIQ METHODOLOGY

3.1 | Overall Framework

We propose DefDIQ, research on the definition, detection, and impact quantification for Android code smell. As shown in Fig. 1, DefDIQ consists of three primary research components: (1) define 15 novel security-related Android code smells and offer targeted repair suggestions, (2) realize DACS to detect targeted smells in Android apps source code automatically, and (3) construct fault counting models to investigate the distribution of custom smells and faults occurrences, quantify the impact of Android code smells and generate repair priorities.

Define Android code smells. Code smells are mainly attributed to deviations from standard coding practices by developers, and are likely leading to security risks and faults. Also, with the specific structure, smells can be identified by relevant information. Previous research findings indicate that Android code smells perform better in security risk prediction²⁹. Therefore, we analyzed the security structure in Android apps and defined 15 novel Android code smells by collecting programming practices, coding specifications, and Android code features related to Android security and defect prediction.

DACS Implementation. Code smells are subjective, so the community lacks the tools to identify the code smells we defined. Accordingly, we design a lightweight automatic detection tool–DACS for the 15 custom smells. DACS implements detection algorithms based on detailed definitions and generates a smell rule library. Firstly, DACS extracts semantics in the leaves and syntax in the non-leaves, effectively representing the semantic and structural information of the source code via AST³⁰. Then search the AST with the detection rule library and identify the code smells. DACS implements a detection subsystem for each code smell and presents the final result as a .csv file or an interface. The source code of DACS is publicly available on GitHub at <https://github.com/strongcat0325/DACS.git>.

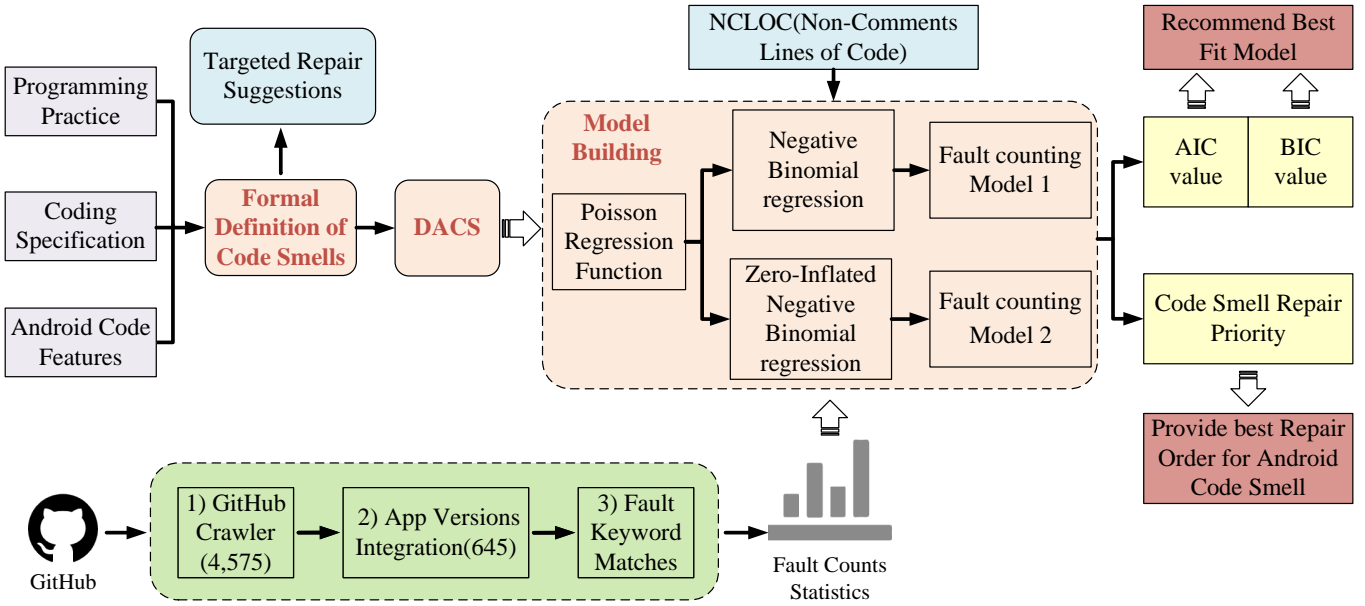


Figure 1 DefDIQ Flowchart

Statistical Analysis. We construct fault counting models with NB and ZINB to investigate the distribution between the occurrences of smells and faults, and statically analyze whether custom smells have an impact on faults. Based on the constructed models, we quantify the impact of code smells and generate repair priority, providing a practical guideline for inexperienced developers. Among them, the fault is an intuitive indicator of app quality. We obtained faults information by the commit specification from GitHub, and gathered faults by version control and keywords match. Then we adopt Android code smells as arguments, while the fault counts as response variables, and apply NB and ZINB regression to construct fault counting models. Ultimately, quantify the code smells' importance based on the parameters in the model and generate repair priorities for all smells by the correlation among smells.

3.2 | Android Code Smell Definition

This paper summarizes and formally defines 15 novel Android code smells and provides targeted repair suggestions. The detailed definition lists are given below.

- **Weak Crypto Algorithm (WCA):** The weak crypto phenomenon usually occurs in Android apps, mainly for two reasons: using weak crypto algorithms; using strong crypto algorithms, but with vulnerabilities in the implementation. Specifically including: 1) Methods with weak crypto hash functions (such as MD2, MD4, MD5, SHA1, or RIPEMD) will trigger code smells; 2) Methods execute the DES algorithm or initialize in AES by ECB mode; 3) Method with RSA algorithm for the key length < 512 bits; 4) Recommend to adopt Cipher.getInstance(RSA/ECB/OAEPWithSHA256AndMGF1Padding) for the crypto algorithm, or hackers will send malicious attacks by received packets.

Repair Suggestions: 1) Recommend developers to adopt SHA256 and SHA3 functions instead of weak crypto hash functions; 2) Recommend developers to adopt AES algorithm, specify CBC or CFB mode available with PKCS5Padding; 3) Recommend developers set the key length to 2,048 bits in the RSA algorithm.

- **Improper Certificate Validation (ICV):** Android provides an embedded process to verify certificates signed by CA. When applying a self-signed certificate, the operating system will verify it through the apps. Unfortunately, developers often fail to implement certificate validation properly. Thus, the communication channel over SSL/TLS is vulnerable to man-in-the-middle attacks: 1) HTTPs communication via WebView, simply handling certificate errors with proceed() in the onReceivedSslError(); 2) Developers customized to implement class X509TrustManager, rewrite method checkClientTrusted() and checkServerTrusted(), but there is no logical code for reviewing the certificate in the method (empty implementation), will generate code smell; 3) Developers adopt a custom implemented class HostnameVerifier, there is no check on the validity of the host name in the verify()

(directly return true); 4) Method with unsafe HostnameVerifier:org.apache.http.conn.ssl.AllowAllHostnameVerifier, org.apache.http.conn.ssl.SSLSocketFactory.ALLOW_ALL_HOSTNAME_VERIFIER from system, the result is equivalent to non-verification.

Repair Suggestions: 1) Recommend developers not to override onReceivedSslError() and not to leave the certificate problem unsolved to avoid the leakage of communication data; 2) For HTTPS websites with certificates signed by authorities, developers can adopt the certificate verification mechanism from Android instead of implementing it; 3) Implement HTTPS domain verification code in HostnameVerifier, i.e., verify whether the domain connected to the HTTPS site and in the SSL certificate is the same; 4) Recommend developers to adopt the secure Android built-in HostnameVerifier: org.apache.http.conn.ssl.SSLSocketFactory.STRICT_HOSTNAME_VERIFIER).

- **Unconstrained InterComponent Communication (UICC):** This refers to the apps reuse the third party components (Activity/ContentProvider/Service/BroadcastReceiver) by setting the attribute exported or IntentFilter in AndroidManifest.xml. For instance, when the exported is true, or the exported value is not set but set to the IntentFilter, the component is exposed to the risk of being invoked by malware.

Repair Suggestions: Recommend developers to set the exported value to false so that external apps are unable to invoke the component. If the developer expects a specific application to visit the component, then exported cannot be set to false, and the attribute permission needs to be set to a custom permission string.

- **Custom Scheme Channel(CSC):** The IntentFilter (with three attributes: action, category, data) in AndroidManifest.xml is for parsing implicit intent. Including: 1) If data android:scheme is specified, then start the intent class to complete communications among the components, and the attackers can visit the APP components by web, causing code smell; 2) If the browser supports Intent Scheme Uri, and sets no filtering rules, then the attacker can visit the browser's files (whether private or public) through JS code, such as stealing cookies files.

Repair Suggestions: 1) Recommend developers not to specify the attribute data android:scheme; 2) If invoking the function Intent.parseUri, then the intent must set strict filtering conditions with at least three modes: addCategory("android.intent.category."), setComponent(null), setSelector(null).

- **Headers Attachment(HA):** When requesting server communication by HTTP, developers usually send data to the server, such as platform number, channel number, system version, and other common information. Such data may be sensitive and rely on header transmission. Specifically, invoking the header file for storing private data is considered a code smell. There are three types of header attachments: set the http header parameters of OkHttpClient, HttpURLConnection and HttpClient.

Repair Suggestions: Before authenticating a third-party service, developers should not store sensitive data in the header.

- **Exposed Clipboard(EC):** In Android apps, the clipboard can store data temporarily in RAM. The app invokes the setText()/getText() or setPrimaryClip()/getPrimaryClip() of the ClipboardManager for writing/reading data to the clipboard. It is insecure to store private data (especially passwords) in the clipboard, which is readable by any app, so we define this behavior as a smell.

Repair Suggestions: Developers should avoid storing data in the clipboard.

- **Broken Web View's Sandbox(BWS):** In the Android SDK, WebView provides a browser allowing users to visit the web in the apps. Among them, the improper utilization of WebView will easily cause arbitrary execution of codes. This code smell occurs for three main reasons: 1) The addJavascriptInterface interface in the WebView generates a Java object, and JS invokes the object's methods to communicate with local APPs. JS can do anything after getting the object, including reading sensitive information on the device's SD card; 2) The Js interface of searchBoxJavaBridge generates JS mapping objects by default, causing the execution of arbitrary code, triggering code smell; 3) Similarly, the Js interface of accessibility and accessibilityTraversal will also lead to the arbitrary execution of codes, generating code smell.

Repair Suggestions: 1) Before Android 4.2, the result to JS is handled by prompt() method; after Android 4.2, Google provides the annotation @JavascriptInterface to avoid vulnerability attacks; 2) Remove the searchBoxJavaBridge interface by invoking the removeJavascriptInterface() method; 3) Remove the accessibility and accessibilityTraversal interfaces.

- **Web View Plain Secret(WPS):** In Android apps, the WebView enables password saving by default: WebView.setSavePassword(true). When the user fills in and agrees to save the password, it will be stored in the database in plain-text and generate code smell.

Repair Suggestions: Suggest the user to close the password saving function: WebView.setSavePassword(false).

- **Web View Domain Not Strict(WDNS):** Android built-in method `setAllowFileAccess()` is for setting whether to allow the WebView to use the file protocol, the default setting is true. When using the file protocol, code smell will be triggered when the following setting is true: 1) `setAllowFileAccessFromFileURLs()`, for setting whether the JS code in the file path is available to visit the local files; 2) `setAllowUniversalAccessFromFileURLs()`, for setting whether the JS code in the file path is available to visit cross-domain sources, such as http domain; 3) `setJavaScriptEnabled()`, for setting whether to allow loading JavaScript.

Repair Suggestions: For general apps, developers should prohibit the file protocol, i. e. `setAllowFileAccess(false)`; for apps requiring special functions with file protocol, developers should set the above attributes to the prohibited state.

- **Data Back Up Any(DBA):** In the Android configuration file, `allowBackup` attribute is for managing the data archive of Android apps. When the value is true, ADB debugging tool can copy and export data without super administrator authentication. Attackers will obtain private data by USB debugging, thus leading to user privacy leakage. We define this behavior as a smell.

Repair Suggestions: Recommend developers set the attribute `allowBackup` in the `AndroidManifest.xml` to false.

- **Global File Readable/Writable(GFRW):** In Android apps, Activity exports data to a file by `openFileOutput(String name, int mode)`. The two parameters represent the name and operation mode of the file, respectively. The details including: 1) When the operation mode is set to `MODE_WORLD_READABLE`, attackers can read sensitive information in the file. This behavior will produce smell; 2) When the operating mode is set to `MODE_WORLD_WRITEABLE`, attackers can write the file's contents freely and destroy the app's integrity, producing code smell.

Repair Suggestions: Developers should identify whether sensitive data is stored in the file. If sensitive data exists, they should: 1) Avoid setting the mode to `MODE_WORLD_READABLE`; 2) Avoid setting the mode to `MODE_WORLD_WRITEABLE`.

- **Configuration File Readable/Writable(CRW):** `SharedPreferences(String name, int mode)` is a common data storage method in Android, which requires two parameters, the name and operation mode of the file. 1) When set mode to `MODE_WORLD_READABLE`, the third-party app can view the file, resulting in leakage of sensitive information; 2) When setting the mode to `MODE_WORLD_WRITEABLE`, the third-party app can delete and change the file, potentially introducing malicious code.

Repair Suggestions: Developers should set the mode of the `getSharedPreferences()` method to `MODE_PRIVATE`.

- **Malicious Unzip(MU):** The Android app will take `ZipInputStream` and `ZipEntry` when decompressing/compressing zip files. The `ZipEntry` provides `getName()` to read the name of the zip file. If the file name has a special string like `../` when decompressing the file, it will produce a code smell.

Repair Suggestions: Recommend developers filter the upper directory string and check the file name when decompressing.

- **SD Visit(SDV):** Android apps store public data in shared external storage by invoking `getExternalStorageDirectory()`. When developers store sensitive information on an external SD card without encryption, code smell will be detected.

Repair Suggestions: Recommend developers to prohibit `getExternalStorageDirectory()`, store sensitive information in the private directory of the application and encrypt the sensitive data.

- **Unsecure Random(UR):** The `SecureRandom` is for obtaining random numbers in Android encryption algorithm. By default, the random number list is generated by `dev/urandom`, and the result is hard to predict. However, when setting the seed, the random number is generated based on a fixed algorithm and seed, so it can be easily predicted. When calling class `SecureRandom`, generating random seeds with the following method is considered a smell: `SecureRandom.SecureRandom(byte[] seed)`, `SecureRandom.setSeed(byte[] seed)` and `SecureRandom.setSeed(long seed)`.

Repair Suggestions: 1) Recommend developers not to generate random numbers by `Random` class; 2) Do not set a seed when using `SecureRandom`.

3.3 | DACS

We developed DACS to detect 15 targeted Android smells in source code. Android code smells are fine-grained method-level code structures and must be identified by analyzing the source codes, so DACS completes the detection based on ASTs from source codes. DACS designs detection rules for each Android code smell and independently realizes the detection sub-module to output the targeted smell results. As depicted in Fig. 2, the detection method mainly contains two modules: (i) convert the APPs source code into ASTs under ESTree standard, and design smell library based on the definitions, then realize the specific

identification algorithms, and (ii) collect the ASTs set and search the AST by specific Android code smell detection algorithm, and count the occurrences of each code smell based on the matching rules.

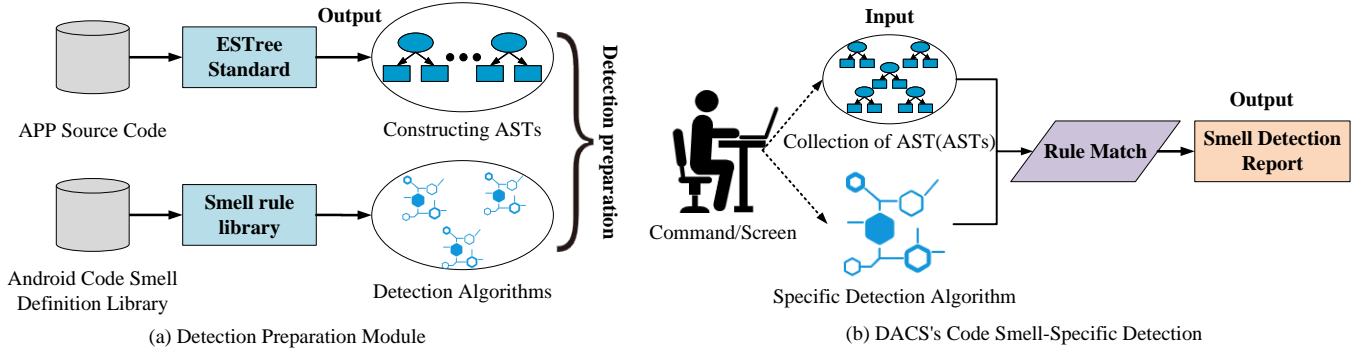


Figure 2 DACS Detection Method

AST facilitates the code smell detection in apps by abstracting the syntactic structure of the source codes in the form of trees, with the relationship nodes representing lexical information and syntactic structure. Then, DACS traverses the AST following the detection algorithm and counts the number of specific code smells in the app. Specifically, (1) depending on the specific definition, summarize the detection keywords; (2) taking the method level as the detection cell, match the keywords within each method body to determine whether the method contains the code smell; (3) calculate the total number of code smells present in the app. Algorithm 1 shows a specific code smell–WCA detection process.

Algorithm 1 WCA Detection Algorithm

```

1: Initialization:  $cnt \leftarrow 0$ 
2: repeat
3:    $method \leftarrow pClassBean.getMethods()$ 
4:    $content \leftarrow method.getTextContent()$ 
5:   if  $content.contrains("MD2" || "MD4" || "MD5" || "SHA-1" || "RIPEMD" || "DES" || "AES/ECB")$  then
6:      $cnt \leftarrow cnt + 1$ 
7:   else if  $content.contrains("RSA")$  then
8:     if  $!content.contrains("RAS/ECB/OAEPWithSHA256 AndMGF1Padding")$  then
9:        $cnt \leftarrow cnt + 1$ 
10:       $macher \leftarrow pattern.matcher(content)$ 
11:      while  $macher.find()$  do
12:        if  $Integer.parseInt(macher.group(2)) < 512$  then
13:           $cnt \leftarrow cnt + 1$ 
14:        end if
15:      end while
16:    end if
17:  end if
18: until  $pClassBean.getMethods()$  is null
Ensure:  $cnt$  //code smell number

```

3.4 | Statistical analysis of impact quantification

DfeDIQ proposes a fault counting model to quantify the relationship between Android code smells and the app's fault count (no association, positive association, negative association, or combined association), providing insights for repair priority. We verify whether code smells and their combinations matter to the fault counts by two-sided hypothesis testing, determine the rejection domain threshold value according to the set significance level, and then determine the model parameters. The discrete regression relationship is established to illustrate which code smells are of no concern to developers and which ones should be

paid special attention to and repaired in time. Among them, the Android code smells are selected from the component *I* custom list; the fault information is extracted from the history submission data of practical open-source Android apps.

3.4.1 | Fault Statistic

The software testing community has different definitions and conceptual descriptions of software problems, and in this paper, the following problems are broadly defined as faults, including (1) fault, describing the incorrect program design in software; (2) vulnerability, representing the static defects in software; (3) error, describing the error results from running defective software; (4) failure, referring to the external software failure behavior observed by users; (5) bug, indicating the general description of software failures and errors; (6) defect, which can cause various bugs. These problems affect quality from different perspectives, including performance, security, memory, energy consumption, portability, and efficiency.

Fault data sources are a critical factor affecting fault counting models. Researchers use data warehouses to collect faults, and these data warehouses can be broadly classified into three categories: private/commercial warehouses, partial public/free warehouses, and public warehouses³¹. Among them, public repositories are adopted by most researchers as a rich and valuable source of fault data (e.g., Bugzilla, JIRA, PROMISE). However, Android apps rely on third-party libraries and update so quickly that no available public data repository has yet emerged. Due to the limitation of traditional manual time-consuming and labor-intensive, researchers mainly collect faults data in large datasets by the method of Zimmermann³².

Android apps on GitHub are diverse, with some in multiple development versions, some in only one version, and others even incomplete. To collect reasonable faults data, we focused the research on the same development version of the Android app and counted the faults by keyword matching. The specific rules for the app version are as follows:

- Incomplete version: search the first complete executable project from the original commit and record the faults between this commit and the latest commit.
- One version: record the faults between this version and the latest one.
- Version number ≥ 2 : record the faults between the latest version and the last committed one.

DefDIQ categorizes the apps by version, locates the commits that contain all fault keywords in the "open" and "closed" problems, and counts the faults. The method is summarized as follows:

- Find commits with the "fail", "vulnerability", "error", "failure", "defect", "bug", "mistake", "fix(ed)" and "update(d)" in the committed and resolved problems. In particular, ignore the case of the keyword while matching.
- Due to the irregularity by the developers, many commits with faulty keywords are not purely faulting repairs. Therefore, it is essential to manually review the commits that successfully match the keywords. Due to massive faults, random sampling detection may be a great choice.

We obtained the statistical information about the app fault through the above steps, and all of the raw fault data is available at <https://github.com/strongcat0325/Fault-Data.git>, as presented in Fig. 3. The original fault counts do not obey the normal distribution with a significant degree of dispersion, and most apps' fault counts are 0.

3.4.2 | Construct fault counting model

By establishing a fault counting model, we can further identify and repair the high-impact defect structures to reduce code review/testing costs. The counting model achieves a linear fitting for discrete data, and the classical discrete counting models are Poisson regression models and zero-inflated Poisson models. Since Poisson distribution has constrained preconditions, the intra-group correlation is not negligible for some interdependent data, so other similar models should be considered.

For determining the suitable counting model, we observed the original data distribution, and the results show that (1) the faults are distributed discretely; (2) the majority of apps do not contain faults, and statistical results show that most mobile apps faults are 0; (3) The faults data set does not obey normal distribution, the variance $>$ the mean, and exists zero-inflated phenomenon. The main reasons for (1)-(2) are that the Android app development cycle is short and evolves rapidly, while the development is not standardized, and testing is insufficient for small and medium-sized projects on GitHub. For (3), we conducted a dispersion test on the faults data by the `dispersiontest()` in R (to verify whether the variance is equal to the mean). The findings indicate that there is a small probability ($p < 0.001$) for the case of data variance = the mean, thus rejecting the original hypothesis. Therefore, NB and ZINB regression models are more suitable for fitting the data than the Poisson regression model.

We design the following scheme to address the discrete data.

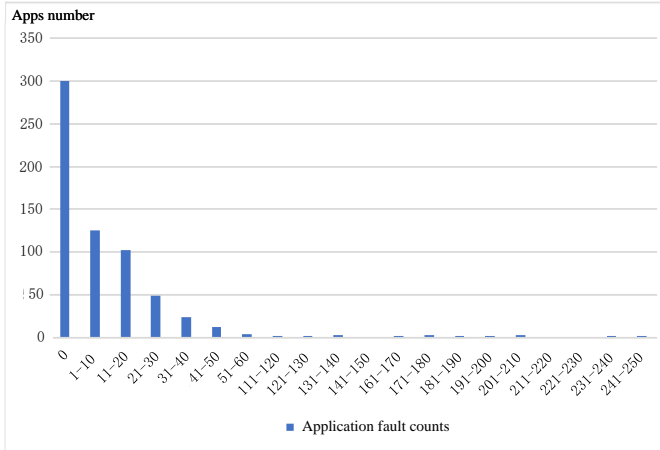


Figure 3 Application Fault Counts Statistics

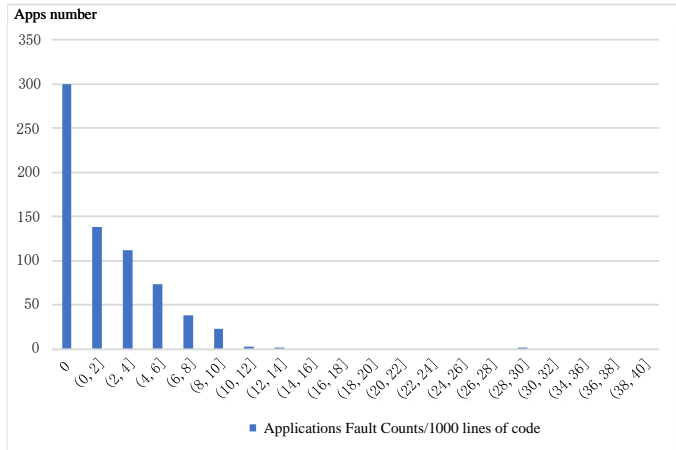


Figure 4 Application Fault Counts Density Statistics

- To minimize the dispersion of fault counts, we take the Density of Faults (i.e., the number of faults per thousand lines of code (NCLOC) in the file) as the response variable, and the measurement unit is faults/KLOC. The processed data are depicted in Fig. 4, but the Density of Faults distribution still does not follow a normal distribution. In the absence of fault data with normal distribution, many counting models failed to work (e.g., Poisson distribution).
- Considering the presence of excessive 0 in the faults data, DefDIQ adopts the NB regression model and ZINB regression model to analyze the relationship between the faults and Android code smells. And the two regression models are inclusive of over-discrete response variables, i.e., allowing the variance > the mean. The statistical results revealed nearly 50% of the apps with 0 faults. Such data are subject to significant bias in the traditional Poisson regression model, while ZINB regression model provide insight to address the excessive apps with "0" fault. In DefDIQ, we apply NB and ZINB to build the counting models and evaluate the model fitness by AIC and BIC metrics. Where the fault counts are dependent variables and some custom Android code smells are selected as arguments.

4 | EXPERIMENTAL SETUP

To investigate and evaluate the performance of DefDIQ in automatically detecting Android code smell and constructing fault counting models, we designed the following research questions (RQs):

- **RQ1:** How effective is DACS detecting custom Android code smells?
- **RQ2:** What is the correlation between Android code smells and the faults in the fault counting model?
- **RQ3:** What is the fitting degree of the fault counting model constructed by NB and ZINB, and to what extent do they affect the quality of the apps?

RQ1 investigates the detection performance of DACS and explores whether the DACS can replace manual detection. Since the code smells detected are novel and customizable, there is no mature tool for comparison. To obtain a reliable reference for the detection results, we invited researchers to accomplish the code smell detection of the apps and take it as a baseline to verify the effectiveness.

RQ2 intends to investigate the distribution between five selected code smells and faults by constructing fault counting models, and observe the correlation between code smells and faults to direct developers and researchers to smells with more significant impacts on security and quality.

RQ3 investigates which regression model is more appropriate for the experimental dataset and quantifies the impact of security code smells on quality based on the best model. Determine the relative importance of code smells on the quality (causing faults), thus assisting in the generation of repair priorities and providing a practical guideline for inexperienced developers.

4.1 | DACS Experiment

In this section, we investigate the DACS performance and verify its effectiveness by measuring the agreement with manual detection results. Specifically, note that, manual detection and DACS only share the same code smell definitions while not sharing any other detection strategies.

4.1.1 | Dataset

The experiment was deployed on 20 open-source Android apps from small & medium-sized projects with high star ratings on GitHub. The 20 apps originate from distinct categories involving learning, Video & Audio, reading, social, game, etc., with different sizes. Table 2 summarises the release and package links associated with the apps.

Table 2 Apps for DACS Effectiveness Evaluation

App ID	Name	Category	Version	App ID	Name	Category	Version
1	NewPipe	Video & Audio	v0.14.2	11	Cupboard	Life	v1.0
2	FxcnBeta	Reading	3.2	12	RITA	Tool	v1.0
3	RGBTool	Tool	v1.4.2	13	Dessarter	learning	v1.0
4	EasyBudget	Life	1.6.2	14	TorchLight	Shopping	v2.4
5	PodTube	Video	1.3.2	15	VINCLÈS	Social	3.12.0
6	Typesetter	Tool	v1.0.1	16	intra42	learning	v0.6.2
7	Snapdroid	Video & Audio	v0.15.0	17	Natibo	learning	v0.2.4
8	LAY	Health	v0.1.3	18	Lexica	Game	v0.11.4
9	Habitat	Tool	v1.0.9	19	Timber	Video & Audio	v1.6
10	WhuHelper	learning	v1.0	20	MateSolver	Game	v2.6

4.1.2 | Experimental Settings

We use DACS to label the custom code smells within Java files in corresponding apps and obtain detection reports. Then we arranged for two experienced Java programmers to manually check the files based on the detailed definitions and count the smells. To minimize the errors caused by individual subjective judgment, the second programmer was arranged to detect again after the first one completed the detection work. The second review focused on whether the first programmer misidentified the smells. The two programmers were independent and without communication during the detection. The results showed only 11 code segments detected by the first programmer were categorized as misidentified by the second. After discussion, six code segments were finally identified as misidentified, and then revised the results. Additionally, to avoid bias, the programmers did not know the experimental details and the specific rules of DACS. By this point, a relatively unbiased and accurate code smell detection result was generated, and we regard it as a comparison baseline.

4.1.3 | Evaluation Method

Code smell is a subjectively defined code structure that possibly negatively influences the app's performance, and there are no open-source, detailed code smell detection methods in the industry. Due to the lack of reliable detection data as a baseline, it is infeasible to adopt precision or recall as evaluation metrics. Consequently, we measure the agreement with manual detection as DACS evaluation metrics. Cohen's Kappa statistic is a common metric for agreement testing, but it is suitable for categorical ratings while having some limitations in count-based models³³. For continuous integer experimental data, Lin's Concordance Correlation Coefficient (CCC) will be better for measuring the agreement of two objects³⁴. Assuming two random variables x and y , the formula for CCC ρ_C is given in Eq. 4.

$$\rho_C = \frac{2\rho\sigma_x\sigma_y}{(\mu_x - \mu_y)^2 + \sigma_x^2 + \sigma_y^2} \quad (4)$$

Where μ_x and μ_y are the means of the two variables, σ_x^2 and σ_y^2 are the corresponding variances, and ρ is the correlation coefficient. When the data set length is N , i.e., $(x_n, y_n), n = 1, \dots, N$. Then the CCC is calculated in Eq. 5.

$$r_c = \frac{2s_{xy}}{s_x^2 + s_y^2 + (\bar{x} - \bar{y})^2} \quad (5)$$

Among them, the mean is obtained by Eq. 6, the variance s_x^2 and covariance s_{xy} are expressed by Eq. 7.

$$\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n \quad (6)$$

$$s_x^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})^2, \quad s_{xy} = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})(y_n - \bar{y}) \quad (7)$$

DefDIQ relied on the CCC evaluation indicator as discrimination for DACS effectiveness, here $N = 15$, representing the 15 custom smells. The data pairs denote each code smell i , the occurrences for manual detection x_i and DACS y_i . Similar to the correlation coefficient, $-1 \leq \rho_C \leq 1$ and $-1 \leq r_C \leq 1$, when the metric approaches +1, it indicates strong agreement between x and y , and as the metric gets closer to -1, it indicates strong disagreement between x and y . Another accepted method is to interpret Lin's CCC metric as Spearman rank correlation coefficient (e.g., < 0.20 denotes poor agreement, while > 0.80 indicates good agreement). In this paper, we follow this interpretation to analyze the agreement strength of the DACS and manual detection results, and the details are shown in Table 1. We rely on the CCC function in the DescTools library in R to obtain Lin's CCC with the confidence interval $1 - \alpha$, ($\alpha = 0.05$).

4.2 | Fault Counting Model Experiment

4.2.1 | Data Collection

In this paper, we collect 4,575 practical Android apps from GitHub and count the faults based on keywords (see Section 3.4.1 for the details). To guarantee the correspondence between the faults and apps, DefDIQ records the faults within one development version, so the 4,575 original apps were reduced to 645. We randomly selected 516 apps as training data and the rest as prediction data. The obtained original data is shown in Fig. 3.

Two main factors were considered in determining the model parameters:

- i The model is vulnerable to over-fitting with excessive parameters or insufficient sample data. Through investigation, we found that most studies constructed models with ≤ 6 parameters^{10,35,36,37}, so we set around 5 parameters for the fault counting model.
- ii The correlation between the parameters. The higher correlation between code smells implies that some smells are redundant and multicollinear, which leads to (1) an unnecessary increase in the feature dimensions, (2) a high tendency to over-fit the model, and (3) interaction between features will reduce model credibility.

To determine suitable parameters to complete the model construction, we detected 15 smells by DACS and counted the occurrences of each smell in 4,575 apps, and the percentages are shown in Fig. 5. Smells with higher occurrences should be selected for model construction to avoid large model errors caused by insufficient data.

Besides, we measured Spearman rank correlation to discover the relationship among 15 custom Android code smells based on DACS detection results. Fig. 6 demonstrates the confusion matrix of Spearman correlation coefficients for code smells, which can visualize the magnitude of the correlation. We found that the correlation among the code smells was quite low (< 0.30), except for the strong correlation for WCA with CSC, MU, SDV, MU with SDV, and HA with MU, SDV. These code smells are independent and will not affect each other, meeting the prerequisites of some regression models while avoiding the experiment's randomness and ensuring the experiment design's rationality.

Finally, we selected WCA, HA, DBA, MU, and SDV as the arguments to fit the faults counting model. In addition, we take the code size into account (NCLOC) as an additional factor to investigate whether file size affects faults.

Fig. 7 shows a matrix scatter matrices of the five smells, provides the faults distribution in the apps, and visualizes the changing relationship between the smells and faults. Furthermore, Fig. 7 illustrates the difference in the density and distribution of code smells in apps, where the horizontal and vertical coordinates are two different code smells.

4.2.2 | Experimental hypothesis

This part aims to investigate whether the five code smells targeted are more likely to be associated with fault-prone files. We adopted a two-sided hypothesis test with a significance level of $\alpha = 0.05$ to determine the model parameters by p-value. P-value

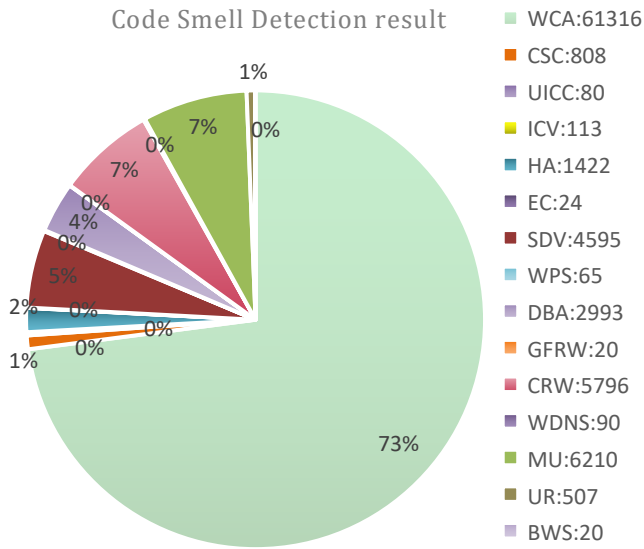


Figure 5 Code Smell Detection Results Statistics

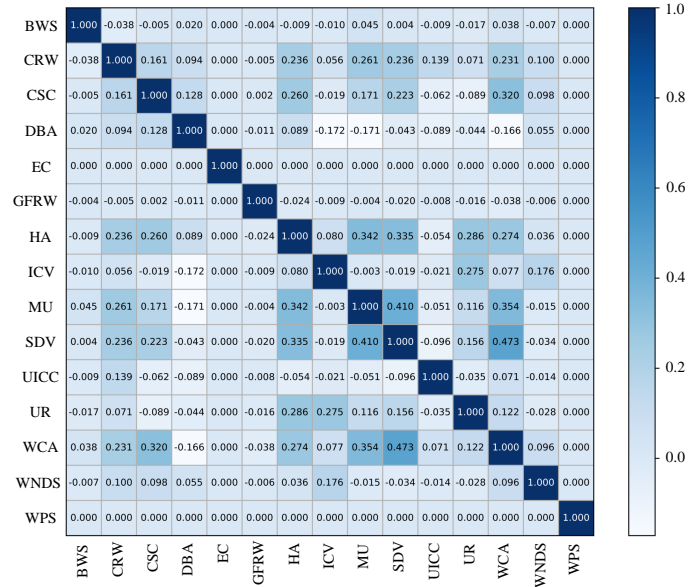


Figure 6 Correlation between Android code smells

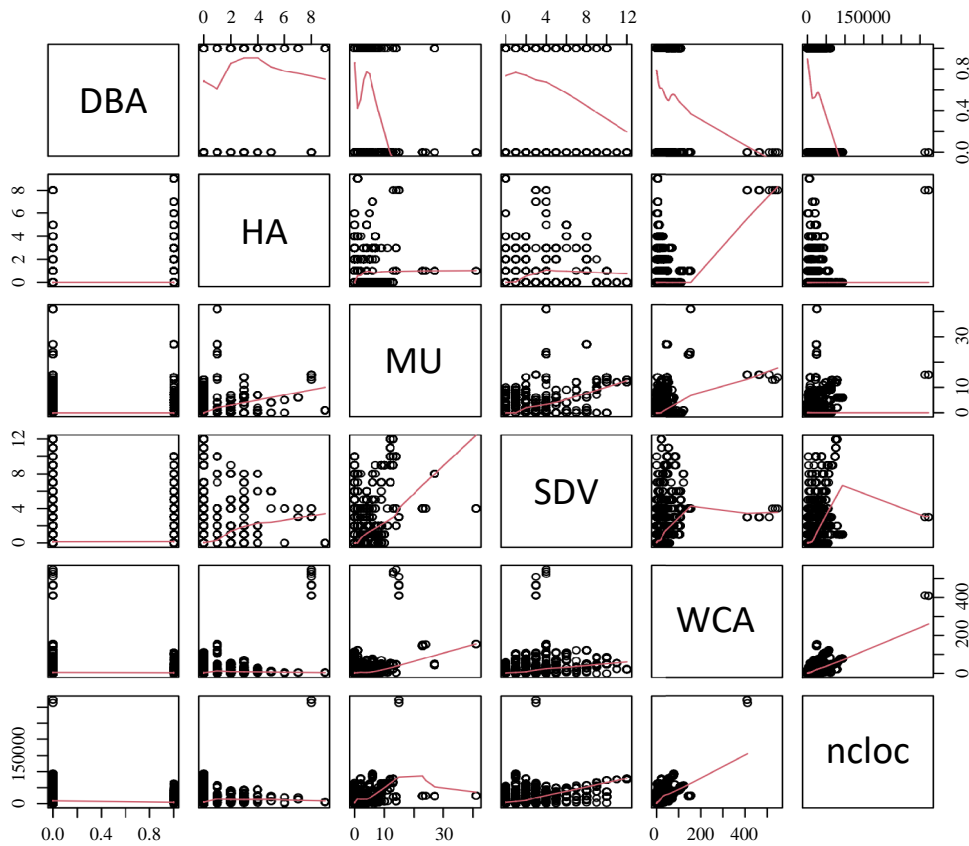


Figure 7 A Scatter Plot Matrix of Code Smells

is an indicator for judging whether to accept the original hypothesis, as explained in Table 3. Specifically, we test the following hypothesis at the significance level of $p < 0.05$:

Hypothesis 1: WCA has no effect on the fault counts in apps (either alone or in combination).

Hypothesis 2: HA has no effect on the faults counts in apps (either alone or in combination).

Hypothesis 3: DBA has no effect on the faults counts in apps (either alone or in combination).

Hypothesis 4: MU has no effect on the faults counts in apps (either alone or in combination).

Hypothesis 5: SDV has no effect on the faults counts in apps (either alone or in combination).

Hypothesis 6: NCLOC has no effect on the faults counts in apps (either alone or in combination).

Table 3 P-value explanation

P-value	Occasionality	Original hypothesis	Statistical Significance
P>0.05	Probability of random error occurring occasionally > 5%	Can not reject the original hypothesis	Code smell has no significant impact on the application fault counts
P<0.05	Probability of random error occurring occasionally < 5%	Can reject the original hypothesis	Code smell has significant impact on the application fault counts
P<0.01	Probability of random error occurring occasionally < 1%	Can reject the original hypothesis	Code smell has a strong significant impact on the application fault counts

4.2.3 | Experimental Settings

DefDIQ constructs a simple and suitable fault counting regression model with minimal parameters under the condition of highly fitted. Hence, we apply the NB and ZINB to model the relationship between fault counts and code smells, respectively. P-values overrode partial hypotheses to filter out some code smells and their combinations affecting faults. After repeated model evaluation and selection, they finally determined the free parameters (code smells). Specifically, the gradual regression steps for model construction are as follows.

- The first-order model contains all independent arguments, i.e., five custom code smells, NCLOC, and combinations. It separates and analyzes the effect of each argument on the faults.
- Select arguments from the first-order model with significant correlation to the fault counts (p-value < 0.05) and build a simpler data model.
- Repeat the previous task, gradually constructing simpler and better-fitting models for the data.

4.2.4 | Evaluation Method

To compare the fitness of the NB and ZINB regression model on the relationship between faults and code smells, we adopted the AIC (Akaike Information Criterion) and BIC (Bayesian Information Criterion) values of the fault counting models under both algorithms to select the counting model with the higher fitting degree. AIC is a weighted function of fitting accuracy and parameters number, and the contribution is to find the accurate model containing the least free parameters based on information entropy. And the AIC compares the deviation of the fitted values of the fault counting model constructed by the regression function with the true values to assess the model's goodness. AIC is obtained by Eq. 8.

$$AIC = -2\ln(L) + 2\alpha \quad (8)$$

Where L denotes the likelihood function, and α represents the number of free parameters that can be estimated. AIC is inversely related to whether the model can better reveal a significant relationship between the response variables and arguments. The smaller the AIC value, the smaller the difference between the fitted and true values of the model, i.e., the higher the fitting degree of the fault counting model.

BIC is an evaluation metric in Bayesian statistics for choosing among two or more alternative models. Although the model selection criteria of AIC and BIC are distinct, there are many similarities, such as the similar interpretation of penalty terms. As shown in Eq. 8 and Eq. 9, when the likelihood value is larger (more accurate model), and the number of free parameters in the model is larger (simpler model), the AIC and BIC value will be smaller, indicating better model performance. For larger data volume, BIC sets a penalty parameter $\ln(n)$ with a larger weight than AIC to determine the optimal parameters and alleviate the overfitting of the model. In addition, BIC also considers the sample number. When there are excessive data samples, the evaluation metric can balance the model in terms of accuracy and complexity.

$$BIC = -2\ln(L) + \alpha\ln(n) \quad (9)$$

Where L denotes the likelihood function, α is the parameter estimated by the model, and n is the sample size (the observations or data points number). The smaller the value of the BIC, the better the model performs.

5 | EXPERIMENTAL RESULTS

5.1 | Answer to RQ1

In this experiment, we took 15 custom Android code smells as detection objects and statistically recorded the code instances of each code smell (with the method level). DACS detected a total of 322 code odor instances, and the most common ones are WCA (163 instances), UR (61 instances), and DBA (17 instances). Statistically, the 20 open-source apps contain 2,232 java files, with 283 files containing the defined code smells.

Table 4 Code Smell Detection Results

Code Smell	WCA	ICA	UCE	ISU	HA	EC	WRCE	WSP	WDNS	DBA	GFRW	PRW	MU	SDV	UR
Manual check result	158	2	1	2	12	2	4	3	6	17	15	8	5	13	57
DACS check result	163	2	2	2	8	3	5	5	8	17	16	11	6	9	61

As shown in Table 4, DACS and manual detection are consistent in detecting ICA, ISU, and DBA, and the results are closer in detecting other code smells. Moreover, we find the number of code smells detected by DACS is mostly higher than the manual detection results. The main reason is that manual detection will combine the context-sensitive code to determine whether the developers deal with some dangerous behaviors. If there is a corresponding handling method, our programmer will not categorize this code structure as a smell, even if the code contains vulnerable fragments. DACS relies purely on matching rules for detection, which is not very flexible, so there is a certain amount of over-reporting and false judgments.

Table 5 Lin's Concordance of Manual and DACS Detection Results

est	lwr.ci	upr.ci
0.9994	0.9976	0.9996

For the 15 custom Android code smells, DACS can identify almost all code smell instances in Android apps correctly. As shown in Table 5, Lin's correlation coefficient r_c reaches 0.9980264 with a confidence interval (99.5%, 99.9%), indicating a high agreement degree for DACS and manual detection results. Only in two cases did the results not meet the desired level, and there will be missed detection, namely malicious unzipping code smell (MU) and header file code smell (HA).

We analyzed the unidentified smell instances by DACS and explored the reasons for the failure to identify them. Eventually, we discovered that the DACS tool missed some smells detection due to the inconsistent compression libraries used for classes affected by the MU and classes considered in the detection rules. E.g., in the VINCLES project (id=15), `com.icetech.silicompressor` package defines class `SiliCompressor` for video, audio, and other files (Gallerypresenter, ChatPresenter, FileUtils) for compression. However, DACS has not yet extended the `SiliCompressor` library, so DACS will not identify such code smells. Meanwhile, the reason for HA smells being ignored to be identified is roughly the same. While Android apps request communication with the server, except for the common request libraries `HttpClient` and `Okhttp` provided in the Smell rule library, there are also higher-level encapsulation libraries such as `Feign`, `Retrofit`, and `Volley`, so DACS can hardly identify some complex smells correctly.

In the meantime, RQ1 provides a potential method to improve the detector's performance. In the future, we will consider adding various libraries to DACS for Android developers. With the support of more third-party libraries, DACS performance will be higher.

Main findings for RQ1: Except for the omission of MU and HA, for most code smells, the detection result of DACS strongly agrees with the manual detection. Lin's concordance correlation coefficient r_c reaches 0.99, demonstrating that DACS can effectively detect 15 custom Android code smells and capture developers' sub-optimal coding practices. Lastly, we can confirm that the DACS is effective in detection accuracy.

5.2 | Answer to RQ2

5.2.1 | Results analysis of the NB fitting model

NB is applied to analyze the relationship between fault counts and code smells. Table 6 presents the first-order interaction model. The *Intercept* row summarizes the average effect of some not explicitly mentioned factors on faults. The remaining rows represent code smells, NCLOC, and combinations associated with the fault counts. Multiplicative interaction terms test the combinations, i.e., if there are two code smells in the test file, then the multiplicative interaction result is 1 (1×1). And the interactive combinations neglect the existence of any other code smells in the files.

Table 6 First-order fault counting model based on NB

	Estimate	Std.Error	z value	$P_r(> z)$
(Intercept)	-0.2615	0.0785	-2.24	0.0138*
NCLOC	0.0013	0.0003	3.58	0.0001***
WCA	-0.5612	0.1843	-2.28	0.0019**
HA	-0.3851	0.3254	-1.08	0.1995
DBA	-0.3021	0.2134	-1.72	0.1123
MU	0.0949	0.0411	2.30	0.0116*
SDV	-0.0996	0.0384	-2.18	0.0199*
NCLOC:WCA	0.0022	0.0013	2.67	0.0035**
NCLOC:HA	0.0007	0.0006	-1.07	0.2810
NCLOC:DBA	-0.0022	0.0020	-1.61	0.1641
NCLOC:MU	-0.0001	0.0002	-2.48	0.0079**
NCLOC:SDV	0.0003	0.0001	2.39	0.0049**
WCA:HA	-0.4969	0.5088	-0.93	0.3321
WCA:DBA	0.2849	0.3899	0.81	0.4392
WCA:MU	0.0193	0.0310	0.49	0.5939
WCA:SDV	-0.0291	0.0864	-0.36	0.6895
HA:DBA	0.8248	0.4201	2.01	0.0498*
HA:MU	0.1811	0.0568	3.09	0.0018**
HA:SDV	-0.1722	0.1298	-1.24	0.2014
DBA:MU	0.0149	0.0309	0.52	0.6278
DBA:SDV	0.0187	0.1019	0.20	0.8455
MU:SDV	-0.0157	0.0101	-1.59	0.1017

As shown in Table 6, the constructed first-order interaction model's residual sum of squared (RSS) over 535 Degrees of Freedom (DF) is 529.8. Since the $RSS < DF$, the model meets the requirements. However, Table 6 shows that NCLOC, WCA, MU, SDV, NCLOC:WCA, NCLOC:MU, NCLOC:SDV, HA:DBA, and HA:MU are the terms significantly associated with the faults ($p\text{-value} < 0.05$). Therefore, we select the above indicators to create a simpler model, shown in Table 7. This model has a better fit compared with the first-order interaction model, with $RSS = 530.1$ and $DF = 547$. NCLOC, WCA, MU, SDV, NCLOC:WCA, and NCLOC:MU significantly correlate with the fault counts ($p\text{-value} < 0.05$).

Table 7 Second-order fault counting model based on NB

	Estimate	Std.Error	z value	$P_r(> z)$
(Intercept)	-0.2499	0.0810	-3.30	0.0021**
NCLOC	0.0019	0.0003	4.68	0.0001***
WCA	-0.5838	0.2047	-2.79	0.0039**
MU	0.0841	0.0219	3.88	0.0001***
SDV	-0.1219	0.0459	-2.58	0.0082**
NCLOC:WCA	0.0019	0.0005	3.48	0.0004**
NCLOC:MU	-0.0003	0.0001	-3.29	0.0011**
NCLOC:SDV	0.0003	0.0002	1.41	0.1579
HA:DBA	-0.0479	0.2588	-0.23	0.8358
HA:MU	0.0272	0.0331	0.84	0.4158

The simplest fault counting model was fitted with the six valid parameters discussed above, and the details are shown in Table 8. The model has $RSS = 526.8$ over $DF = 550$. There is a significant association ($p\text{-value} < 0.05$) between the faults and six code smells, i.e., a significant fault-proneness to the app quality.

Table 8 Third-order fault counting model based on NB

	Estimate	Std.Error	z value	$P_r(> z)$
(Intercept)	-0.2958	0.0759	-3.92	0.0001***
NCLOC	0.0028	0.0003	6.31	0.0000***
WCA	-0.5993	0.2065	-2.99	0.0035**
MU	0.0677	0.0198	3.31	0.0010***
SDV	-0.0806	0.0321	-2.66	0.0121**
NCLOC:WCA	0.0020	0.0006	3.49	0.0003**
NCLOC:MU	-0.0002	0.0001	-2.89	0.0031**

5.2.2 | Results analysis of the ZINB fitting model

We also apply the ZINB model to investigate the distribution relationship between faults and code smells, and the modeling process is similar to NB. After obtaining the significantly correlated parameter individuals, we continue constructing the fault counting regression model until all the parameters in the model are significantly correlated. Instead of listing all detailed model information, only the final model results are presented (see Table 9).

Table 9 Fault counting model based on ZINB

	Estimate	Std.Error	z value	$P_r(> z)$
(Intercept)	-0.2768	0.0821	-4.15	0.0001***
NCLOC	0.0031	0.0002	6.25	0.0000***
WCA	-0.5980	0.2001	-2.75	0.0025**
MU	0.0665	0.0172	3.28	0.0008***
SDV	-0.0146	0.0100	-4.83	0.0031**
NCLOC:WCA	0.0019	0.0005	3.38	0.0002**
NCLOC:MU	-0.0001	0.0001	-2.78	0.0028**

As displayed in Table 8 and 9, the results remain consistent, showing a significant correlation between six code smells and the faults in the apps. However, this correlation is not always positive, some code smells (WCA, SDV, and NCLOC:MU) are negatively correlated with the fault counts. The remaining code smells (NCLOC, MU, and NCLOC:WCA) positively correlate with the fault counts. The results show that several code smell combinations increase the fault-proneness, while some code smells decrease the faults. The reasons for such a phenomenon do not exclude insufficient experimental data, as most projects on GitHub are not oriented to the market, and most of these projects are immature for not following the completed software lifecycle steps. The development is not standardized with insufficient testing and maintenance, and thus fault information and smell data collected are scarce. Developers need to repair and refactor the code smells with positive fault correlation in time and carefully refactor the code smells with negative fault correlation.

Main findings for RQ2: The two-sided hypothesis test eliminated the code smells and combinations with no significant effect on the faults at the significance level $\alpha = 0.05$, which ultimately determined six free arguments in the model construction. The research revealed the following findings among the five code smells, (1) distinct code smells and combinations have different importance, some code smells do indicate fault-proneness, such as NCLOC, MU, and NCLOC:WCA; (2) Some smells are not directly related to faults, e.g., HA and DBA had no significant effect on the fault-proneness.

5.3 | Answer to RQ3

We evaluate the fitness for fault counts and code smells by AIC and BIC values of the fault counting model under the same dataset, and judge the performance of the regression models constructed by NB and ZINB. When the model becomes more complex (k rises, i.e., the feature number grows), the likelihood function also increases, thus reducing the AIC and BIC values, but it may cause overfitting issues when the data is too heavy. In the equations $AIC = 2k - 2\ln(L)$, $BIC = \ln(n) * k - 2\ln(L)$, L represents the maximum likelihood, k refers to the free parameters for estimation (6), and n is the experimental data volume

(516). The AIC and BIC values of the final NB fault counting model are 552.68 and 557.48, respectively, while for the ZINB model, the values are 517.32 and 522.12, obviously indicating that the ZINB model fits the experimental data better than the NB (both AIC_ZINB and BIC_ZINB values are lower than NB). As a fault-based regression, the fault counting model constructed by ZINB performed better than the one based on NB.

For a more intuitive view of the extent to which code smell affects fault counts, we construct the repair priority level according to the best-fitting fault counting model on ZINB. Specifically, considering the direct correlation of the coefficients of the arguments with the faults in the model, we take them as the influence factors of the arguments and map the coefficients to the interval [0,1] by the Sigmoid function. The Eq. 10 for the S function is as follows:

$$f(x) = \frac{1}{1 + e^{-x}} \quad (10)$$

According to the ZINB fault counting model, we obtained coefficients representing the importance levels of three independent code smells (WCA, MU, SDV), and then calculated the importance value by S-function. The importance values for the remaining smells were derived by the Spearman correlation coefficient, where correlations below 0.1 were considered as no correlation. Therefore, we calculated the importance-value of HA, CSC, CRW, UR, DBA, ICV, UICC, and WNDS based on the spearman part > 0.1. Importance values for EC, WPS, BWS, and GFRW could not be obtained due to a lack of research data and low correlation, assigning them the lowest repair level considering their occurrence frequency in the files. The repair priority levels of the 15 code smells are shown in Table 10.

Table 10 Repair priority level on ZINB fault counting model

Code smell	Importance-value	Repair priority level
MU	0.5116	I
CRW	0.5069	II
SDV	0.4964	III
HA	0.4405	IV
CSC	0.4273	V
WCA	0.3548	VI
UR	0.3065	VII
ICV	0.0842	VIII
UICC	0.0706	IX
DBA	0.0549	X
WNDS	0.0148	XI
EC,WPS,BWS,GFRW	empty	XII

Main findings for RQ3: Among the two fault counting models constructed by NB and ZINB regression algorithms, the ZINB-based fault counting model is more suitable for the dataset under work (AIC=517.32, BIC=522.12). We generate fault repair priority levels based on the ZINB fault counting model by quantifying the code smells impact on the faults to assist developers and researchers in focusing on fault-prone code smells. The priority list drives the building of high-quality and low-debt apps by refactoring poor designs with high risk first, thereby increasing the repair benefits.

6 | THREATS TO VALIDITY

Threats to Internal Validity. Code smell is an abstract definition of poor design by researchers from different practice results, and the definition is subjective and without strict criteria. Since the opinion of Android code smells varies among practitioners, most complete the definitions depending on personal practices. To achieve broad acceptance, we carefully investigate technical blogs and conferences related to Android security and code specification, provide the Code Smell Definition, and targeted repair suggestions based on programming practices.

Threats to External Validity. The experimental apps selected for this study are from GitHub, and the validity of the experimental findings for other communities is debatable. Android apps depend on third-party libraries and are developed on a small scale with rapid iterations, and there is no public database. When obtaining fault data, irregular commits by developers lead to the less reliable classification of commits with matching fault keywords as faults. Besides, some projects on GitHub are immature, and most do not follow the complete software lifecycle steps. The project development is not standardized with inadequate

testing and maintenance, and thus fault information and smell data collected are scarce. And the limited dataset under test resulted in a single-digit number of identified code smells. However, given that our apps are from different domains and sizes, we believe the findings still provide valuable information for code smell and app quality studies.

7 | RELATED WORKS

7.1 | Code Smell Research and Detection

Code smell detection is a technique for analyzing and understanding the defined poor design in apps by extracting code metrics from the source code to explore fault-prone design issues. Thus, detection methods, as well as tools, have been continuously adopted to detect code smells in apps^{38,39,40}. Traditional Java static code metric detection is unable to precisely focus on poor design and implementation in Android apps for Android features. With the maturity of Android technology, researchers started to put more effort into detecting Android code smells. Huang et al.⁴¹ proposed four mapping metrics and implemented a static detection tool, HBSniff, to detect 14 code smells in object-oriented programming. Mario et al.⁴² implemented DECOR to detect several object-oriented code smells in apps. Hecht et al. developed PAPRIKA to identify four Android-specific code smells based on manually defined detection rules⁴³. Palomba et al.⁴⁴ implemented aDoctor to detect 15 Android-specific code smells with the abstract syntax of the source code. Fatima et al. extended custom rules based on the AL to detect and correct 12 Android code smells, which ultimately outperformed PAPRIKA⁴⁵. Ghafari et al. investigated 28 Android code smells indicating security and discussed relevant elimination methods².

7.2 | Effects of Smell on Faults

As an indicator of app quality, the fault counts can help quantify the influence implied by Android code smell. The fault prediction models have been researched to locate faulty code, thereby improving software quality and using resources better. Hall et al.⁴⁶ analyzed 36 fault-related works systematically and found that source code-based spam filtering techniques performed well in prediction^{47,48}. Static code metrics have also contributed to fault prediction, but the predictive performance is relatively poor individually. Zhou et al. applied various complexity metrics to logistic regression models and found that NCLOC outperformed other metrics⁴⁹. Furthermore, fault prediction models using only the NCLOC perform better than other source code metrics, and LOC positively correlates with fault counts. Ostrand et al.⁵⁰ and Hongyu et al.⁵¹ both reported that NCLOC is a generic, simple and effective fault prediction metric. Although some researches indicate that the predictive performance of NCLOC is inferior to other metrics, overall, NCLOC has significant advantages in fault prediction models^{52,53}.

8 | CONCLUSIONS

Software suffers from various security and quality threats within the continuous development of attack and counter-attack techniques. The security of Android apps with a high market share is even more critical because of the vast amount of users' sensitive data involved. To reduce the vulnerable code practice, we promote security design practices. In this paper, we customize 15 novel Android code smells and provide targeted suggestions for eliminating or mitigating them. Then we developed a lightweight tool to complete automatic detection and find their universality. We investigated the distribution between code smells and faults by constructing fault counting models, then generated repair priority levels based on the experimental findings. The majority of the app datasets collected contained at least four security smells, and we believe that the identified code smells are a valuable indicator of security and quality issues. The experimental results show that code smells do have an impact on the faults in apps, and the code smell repair priority generated can provide a practical baseline for researchers and inexperienced developers.

In the future research works, we are committed to a larger study involving DACS performance repair, a more extensive Android apps faults repository, and the validation of repair suggestions, realizing more mature research on the topic.

ACKNOWLEDGMENTS

This work is supported in part by National natural science foundation of China (Key Program): Software testing technology for security-critical deep learning system(61832009) and National natural science foundation of China (Key Program): Data-driven testing methodologies and echnologies for intelligent software systems(61932012).

References

1. Lewowski T, Madeyski L. How far are we from reproducible research on code smell detection? A systematic literature review. *Information and Software Technology* 2022; 144: 106783.
2. Ghafari M, Gadiant P, Nierstrasz O. Security smells in android. In: IEEE. ; 2017: 121–130.
3. StatCounter (2022) Mobile operating system market share worldwide. Accessed 6 October 2022: <https://gs.statcounter.com/os-market-share/mobile/worldwide>.
4. Senanayake J, Kalutarage H, Al-Kadri MO, Petrovski A, Piras L. Developing Secured Android Applications by Mitigating Code Vulnerabilities with Machine Learning. In: ASIA CCS '22. The organization. Association for Computing Machinery; 2022: 1255–1257.
5. Alkandari MA, Kelkawi A, Elish MO. An Empirical Investigation on the Effect of Code Smells on Resource Usage of Android Mobile Applications. *IEEE Access* 2021; 9: 61853–61863.
6. Boutaib S, Bechikh S, Palomba F, Elarbi M, Makhoul M, Said LB. Code smell detection and identification in imbalanced environments. *Expert Systems with Applications* 2021; 166: 114076.
7. Fowler M. Refactoring: Improving the Design of Existing Code. In: Lecture Notes in Computer Science. The organization. ; 2002: 256–256.
8. Brown WH, Malveau RC, McCormick HWS, Mowbray TJ. *AntiPatterns: refactoring software, architectures, and projects in crisis*. John Wiley & Sons, Inc. . 1998.
9. Habchi S, Moha N, Rouvoy R. Android code smells: From introduction to refactoring. *Journal of Systems and Software* 2021; 177: 110964.
10. Hall T, Zhang M, Bowes D, Sun Y. Some code smells have a significant but small effect on faults. *TOSEM* 2014; 23(4): 1–39.
11. Paulo Sobrinho dEV, De Lucia A, Almeida Maia dM. A systematic literature review on bad smells–5 w's: which, when, what, who, where. *IEEE Transactions on Software Engineering* 2018; 47(1): 17–66.
12. Walter B, Alkhaeir T. The relationship between design patterns and code smells: An exploratory study. *Information and Software Technology* 2016; 74: 127–142.
13. Fenske W, Schulze S. Code smells revisited: A variability perspective. In: VaMoS '15. The organization. Association for Computing Machinery; 2015; New York, NY, USA: 3–10.
14. Xian Zhang JZ. A slice-granularity defect prediction method based on code naturality. *Ruan Jian Xue Bao/Journal of Software* 2021; 32(07): 2219–2241.
15. Viega J, McGraw G, Mutdosch T, Felten EW. Statically scanning java code: Finding security vulnerabilities. *IEEE software* 2000; 17(5): 68–74.
16. Lu Z, Mukhopadhyay S. Model-based static source code analysis of java programs with applications to android security. In: 2012 IEEE 36th Annual Computer Software and Applications Conference. The organization. ; 2012: 322–327.
17. Mannan UA, Ahmed I, Almurshed RAM, Dig D, Jensen C. Understanding Code Smells in Android Applications. In: MOBILESoft '16. The organization. ; 2016; New York, NY, USA: 225–234.
18. Hecht G, Rouvoy R, Moha N, Duchien L. Detecting antipatterns in android apps. In: 2nd ACM International Conference on Mobile Software Engineering and Systems. The organization. ; 2015: 148–149.
19. Jan Reimann UA. A Tool-Supported Quality Smell Catalogue For Android Developers. *Softwaretechnik-Trends* 2014; 34(2).
20. Hecht G, Moha N, Rouvoy R. An empirical study of the performance impacts of android code smells. In: MOBILESoft '16. The organization. ; 2016; New York, NY, USA: 59–69.

21. Gupta PL, Gupta RC, Tripathi RC. Analysis of zero-adjusted count data. *Computational Statistics & Data Analysis* 1996; 23(2): 207–218.
22. Greenwood M, Yule GU. An inquiry into the nature of frequency distributions representative of multiple happenings with particular reference to the occurrence of multiple attacks of disease or of repeated accidents. *Journal of the Royal statistical society* 1920; 83(2): 255–279.
23. Everttt B. Distributions in Statistics: Discrete Distributions. 1970: 482–483.
24. Lambert D. Zero-inflated Poisson regression, with an application to defects in manufacturing. *Technometrics* 1992; 34(1): 1–14.
25. Greene WH. Accounting for excess zeros and sample selection in Poisson and negative binomial regression models. 1994.
26. Caudill SB. Estimating the circle closest to a set of points by maximum likelihood using the BHHH algorithm. *European Journal of Operational Research* 2006; 172(1): 120–126.
27. Hopkins WG. *A new view of statistics*. Will G. Hopkins . 1997.
28. Fieller EC, Hartley HO, Pearson ES. Tests for rank correlation coefficients. I. *Biometrika* 1957; 44(3/4): 470–481.
29. Gong A, Zhong Y, Zou W, Shi Y, Fang C. Incorporating Android Code Smells into Java Static Code Metrics for Security Risk Prediction of Android Applications. In: The organization. IEEE; 2020: 30–40.
30. Fang C, Liu Z, Shi Y, Huang J, Shi Q. Functional Code Clone Detection with Syntax and Semantics Fusion Learning. In: ISSTA 2020. The organization. ; 2020; New York, NY, USA: 516–527.
31. Rathore SS, Kumar S. A study on software fault prediction techniques. *Artificial Intelligence Review* 2019; 51(2): 255–327.
32. Kim S, Zimmermann T, Whitehead EJ, Zeller A. Predicting Faults from Cached History. In: ISEC '08. The organization. ; 2008; New York, NY, USA: 15–16.
33. Fleiss JL, Cohen J, Everitt BS. Large sample standard errors of kappa and weighted kappa.. *Psychological bulletin* 1969; 72(5): 323.
34. Lawrence I, Lin K. A concordance correlation coefficient to evaluate reproducibility. *Biometrics* 1989: 255–268.
35. Muse BA, Rahman MM, Nagy C, Cleve A, Khomh F, Antoniol G. On the Prevalence, Impact, and Evolution of SQL Code Smells in Data-Intensive Systems. *CoRR* 2022; abs/2201.02215.
36. Li W, Shatnawi R. An empirical study of the bad smells and class error probability in the post-release object-oriented system evolution. *Journal of Systems and Software* 2007; 80(7): 1120–1128.
37. Olbrich S, Cruzes DS, Basili V, Zazworka N. The evolution and impact of code smells: A case study of two open source systems. In: 3rd International Symposium on Empirical Software Engineering and Measurement. The organization. ; 2009: 390–400.
38. Moha N, Guéhéneuc YG, Duchien L, Le Meur AF. Decor: A method for the specification and detection of code and design smells. *IEEE Transactions on Software Engineering* 2009; 36(1): 20–36.
39. Patnaik A, Padhy N. A Hybrid Approach to Identify Code Smell Using Machine Learning Algorithms. *International Journal of Open Source Software and Processes (IJOSSP)* 2021; 12(2): 21–35.
40. Sharma T, Efstathiou V, Louridas P, Spinellis D. Code smell detection by deep direct-learning and transfer-learning. *Journal of Systems and Software* 2021; 176: 110936.
41. Huang Z, Shao Z, Fan G, Yu H, Yang K, Zhou Z. HBSniff: A static analysis tool for Java Hibernate object-relational mapping code smell detection. *Science of Computer Programming* 2022; 217: 102778.

42. Linares-Vásquez M, Klock S, McMillan C, Sabané A, Poshypanyk D, Guéhéneuc YG. Domain Matters: Bringing Further Evidence of the Relationships among Anti-Patterns, Application Domains, and Quality-Related Metrics in Java Mobile Apps. In: ICPC 2014. The organization. ; 2014; New York, NY, USA: 232–243.
43. Hecht G, Benomar O, Rouvoy R, Moha N, Duchien L. Tracking the software quality of android applications along their evolution (t). In: 30th IEEE/ACM International Conference on Automated Software Engineering (ASE). The organization. ; 2015: 236–247.
44. Palomba F, Di Nucci D, Panichella A, Zaidman A, De Lucia A. Lightweight detection of android-specific code smells: The adocor project. In: 24th International Conference on Software Analysis, Evolution and Reengineering (SANER). The organization. ; 2017: 487–491.
45. Fatima I, Anwar H, Pfahl D, Qamar U. Detection and Correction of Android-specific Code Smells and Energy Bugs: An Android Lint Extension. In: . 2767 of *CEUR Workshop Proceedings*. The organization. ; 2020: 71–78.
46. Hall T, Beecham S, Bowes D, Gray D, Counsell S. A systematic literature review on fault prediction performance in software engineering. *IEEE Transactions on Software Engineering* 2011; 38(6): 1276–1304.
47. Mizuno O, Ikami S, Nakaichi S, Kikuno T. Spam filter based approach for finding fault-prone software modules. In: MSR'07: ICSE Workshops. The organization. ; 2007: 4.
48. Mizuno O, Kikuno T. Training on errors experiment to detect fault-prone software modules by spam filter. In: ESEC-FSE '07. The organization. ; 2007; New York, NY, USA: 405–414.
49. Zhou Y, Xu B, Leung H. On the ability of complexity metrics to predict fault-prone classes in object-oriented systems. *Journal of Systems and Software* 2010; 83(4): 660–674.
50. Ostrand TJ, Weyuker EJ, Bell RM. Where the bugs are. *ACM SIGSOFT software engineering notes* 2004; 29(4): 86–96.
51. Zhang H. An investigation of the relationships between lines of code and defects. In: IEEE International Conference on Software Maintenance. The organization. ; 2009: 274–283.
52. Bell RM, Ostrand TJ, Weyuker EJ. Looking for bugs in all the right places. In: ISSTA '06. The organization. ; 2006; New York, NY, USA: 61–72.
53. Olague HM, Etzkorn LH, Messimer SL, Delugach HS. An empirical validation of object-oriented class complexity metrics and their ability to predict error-prone classes in highly iterative, or agile, software: a case study. *Journal of software maintenance and evolution: Research and practice* 2008; 20(3): 171–197.

How to cite this article: Yi Z., MY. Shi, JW. He, CR. Fang, and ZY. Chen (2022), Security-based code smell definition, detection, and impact quantification in Android, *Q.J.R. Meteorol. Soc.*, 2022;00:1–6.