

## DroidSD: An Efficient Indexed Based Android Applications Similarity Detection Tool

JUNAID AKRAM, ZHENDONG SHI, MAJID MUMTAZ AND PING LUO

*Key Laboratory of Information System Security*

*School of Software*

*Tsinghua University*

*Beijing, 100084 P. R. China*

*E-mail: {znd15; szd15; maji16}@mails.tsinghua.edu.cn; luop@mail.tsinghua.edu.cn*

Android is becoming more and more popular in recent years. Meanwhile, it has been noticed that the security threats are also increasing with the passage of time because of the code reuse from other applications. In this paper, we propose and design DroidSD, a clone detection tool for Android applications, that helps to detect different types of code clones from APK's source code. A prototype has been developed and implemented to detect clones in Android applications. DroidSD detects Type-1, Type-2 and Type-3 (near-miss) clones in Android applications at the source code level with high accuracy rate, which was not possible in previous Android similarity detection techniques. DroidSD can detect full and partial level similarity between applications. We evaluate DroidSD on real time data-set and count the Recall and Precision on BigCloneBench, which is quite significant. The real time data-set includes 30,500 Android applications including top one, *i.e.* Chrome, Firefox, Gmail, WhatsApp, GoogleMap, Google-PlayStore, Baidu and BigCloneBench.

**Keywords:** clone detection, maintaining APK code, Android apps re-usability, plagiarism detection, Apps similarity detection, information security

### 1. INTRODUCTION

Recently, smart phones are incredibly popular and very widely used in modern life. Android is one of the main smart phone operating system used globally. As we know that Android is an open source operating system for mobile phones, that's why it has been supported by many well-featured applications. There are bundle of Android app stores, which are providing the facility to download latest and updated versions of Android applications<sup>1</sup>. For example, there are more than 3.8 million applications in Google play store [1] and more than 50 billion downloads [2]. These apps are actually providing very useful features in online payment systems, but meanwhile becoming the target for criminals to fraud. Mobile browsers also relied on sensitive security operations, such as online payments and transactions [3-5]. Android is very easy fraud target for attackers because anyone can make their own apps and upload on official app stores. Even they can download the original APK files and by using reverse engineering techniques, they can decompile the source code and rebuild app after making some changes in it [6, 7].

Because of very fast demand of new features in Android applications, the developers are copying code from other apps and tries to reuse these fragments by pasting in other source code sections with or without modifying the code, this type of reuse source code

---

Received April 7, 2018; revised July 1 & August 27, 2018; accepted August 29, 2018.

Communicated by Shiao-Li Tsao.

<sup>1</sup><http://www.businessofapps.com/guide/app-stores-list/>

approach is called code cloning, and the pasted code called as cloned code of the original code. It is very adapting approach, especially in application development activities. However, during or after development process it is quite difficult to say which code fragment is the original and which was copied. Clones in source code actually bring a big trouble in applications security and maintenance [7]. The previous research work shows that a significant amount of 7%-23% of source code was actually cloned in large systems [8]. A lot of research have been done in the field of open source code clone detection in large scale systems, but clone detection in Android applications is still an open chapter. There are some techniques to check the similarity between Android application [9-11], but these techniques check the similarity on basis of it's user interface, application signatures, application versions and uploaded patterns. But in our approach we avoid these comparing factors to check the similarity in applications and focused on analysis the source code of applications. Many researchers have explored the techniques to detect malware in Android applications [2, 12]. Alam [13] and Chen [2] use clone detection techniques to detect malware in APK. Very rare researchers detect code clones at source code level by decompiling Android applications and evaluate the source code. In our approach, we have detected code clones at very deep level of source code by extracting the main feature of APK source files.

To make sure the security and reliability of Android applications, we develop DroidSD to detect the apps similarity, copy paste code from APK source files and injected code fragments, which can be a malware code. Meanwhile it can detect both partial application similarity and full application similarity. We have experimented DroidSD solution on the source code of 30,500 Android applications, we get source code of these apps by using reverse engineering techniques. It detects Type-1, Type-2, Type-3 code clones from Android applications and retrieve the results in the form of similar code fragments. DroidSD worth more for application markets rather than an approach embedded in Android devices or used by end users. The particular apps we actually are interested in finding code clones are those that copy code from other apps, or repackage existing apps. DroidSD is a semantic based clone detection approach, which is scalable, incrementable and can be extended to large scale Android source repositories.

### 1.1 Types of Clones

**Type-1 (Exact clones):** The identical code fragments, which are the exact copies of code, except blanks, layouts, whitespaces and comments [14].

**Type-2 (Parameterized/Renamed):** Two Syntactically identical code fragments are similar except for variations in literals, names of variables, types, and functions [14].

**Type-3 (Near miss clones/Gapped clones):** Two copied code fragments with further modifications such as added or removed statements, the use of different literals, identifiers, types, whitespaces, comments and layouts [14].

**Type-4 (Semantic clones):** Two code fragments that perform the same computation but implemented by different syntactic variants. Or they are semantically similar, without being syntactically similar [14].

The cloned code is actually the similar source code fragment between two applications, where cloned type describes the degree of similarity between code fragments. Type-1 clones are the totally similar code fragments with 100% similarity because the code is fully identical. Type-2 have the similar source code with 90% similarity because in Type-2 clones there are little variations in literals, names of variables, types, and functions but the source code is also similar. In case of Type-3 clones, where new lines added in the original source code or maybe deleted, it is still suspected code which shows the

similarity of 60% to 80% between source code fragments, which have been further categorized and explained in case study of Section 6. So, the similarity between Android source code varies with the detected code clone types.

## 2. RELATED WORK

Detecting and evaluating code clones in applications or code bases have been very useful to many software engineering techniques such as re-factoring and bug detection [15]. Researchers have explored and proposed different ways to identify and detect similarity between Android applications [9, 10], similarity between different documents [16] and in open source files [14, 17-21]. Their results and methods have been employed by code clone management tools for open source projects [22]. Prior research has shown that there are many cloned and fake applications in Android markets [9, 23]. Mostly signature based techniques have been used for malware detection in source code [2]. DroidMOSS [23] detect the similarity of two applications on the basis of fuzzy hashing. Jian Chen [2] collect the Android applications which were known to be malware for malware detection. Sungmin Kim [24] proposed a method to detect illegally copied Android application on the network. He extended data objects which were being transmitted from the network through sniffing, assembling and analyzing the packets. DroidClone [13] exposes the code clones in Android applications by using MAIL (Malware Analysis Intermediate Language) technique. This technique uses a specific flow pattern to reduce the obfuscations effect. AnDarwin [6] uses PDG (Program Dependence Graph) approach to detect the similarity of two applications. It only analyzes the application at the level of Java byte code. But in our method we check and identify the similarity at source code level, we use our own developed method to detect code clone fragments in two different applications. Svajlenko [25] have proposed an efficient way to detect large-scale near-miss clone in large scale systems. Recently SourcererCC tool performs code clone detection in big code and it's extended version SourcererCC-1 is an Eclipse plug in, which is using SourcererCC to detect, identify and navigate all clones during software development.

## 3. PROBLEM AND PROPOSED SOLUTION

It is very easy to apply reverse engineering techniques on Android APK files and re-build them by adding some malicious code fragments inside. Mostly previous techniques have focused on similarity detection between application by consider it's user interface, uploaded information, signatures and patterns. But there are very less techniques, which checks the applications similarity at source code level. Even some developers copy some code from other applications and paste into their own application without testing that code, this code fragment can be the malicious code of original applications. So if we can detect that malicious code fragments from APK files we can use that code fragment as malware pattern to check in other applications. Hence, to perform that task code clone detection technique can be best solution. To make sure the security level and reliability of Android APK files, we propose a scalable solution in the form of DroidSD for clone detection in APK source code files. Our purpose is to build a tool, where we can detect the cloned Android applications. DroidSD has been tested on almost 30,500 Android applications, those downloaded from AppChina market. It detects similarity between applications at different granularity. Our proposed solution is defined by two main steps:

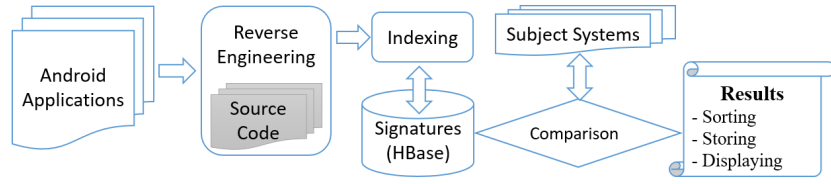


Fig. 1. Top level view of whole system.

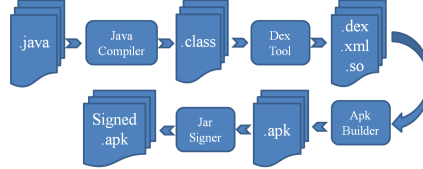


Fig. 2. Building structure of Android applications.

**Step 1: (APK Source Code Collection):** All 30,500 applications decompiled to get DEX files by using reverse engineering techniques. From DEX files we retrieved JAR files, then JAR files further decompiled to a Java source files as shown in Fig. 3.

**Step 2: (APK Code Clone Detection):** We developed DroidSD, a code clone detection tool to detect the same code fragments in different applications to check the similarity between different applications. Our code clone detection approach is hybrid (semantic), which detects Type-1, Type-2, Type-3 code clones in Android applications.

#### 4. APK CODE COLLECTION AND DECOMPILED

Google Play Store is the largest distributed Android channel for Android applications, which is actually offering global coverage for the big audience [7]. But we used AppChina<sup>2</sup> platform to download APK files. The reason we use AppChina store is that it has the application's review criteria [2]. Every application would be checked from developer's side before it released. AppChina has more than 30 million users and about 600 million applications have been downloaded every month<sup>3</sup>.

##### 4.1 APK Crawler

For the collection of APK files, we wrote a web crawler, which was downloading the top ranking applications from market including basic information, *i.e.* Name, Path, Size, Downloaded Frequency and Category).

##### 4.2 Reverse Engineering (.apk to .java source)

We used reverse engineering techniques to get the source code from APK files. The development process of an Android app is shown in Fig. 2. Our most concern files in apk is DEX files because these are the files holding all source code of Java class files. Our performed reverse engineering steps have been shown in Fig. 3.

Reverse Engineering on 30,500 Android applications was done in first part through the decompilation of all DEX files to Java source code files. We wrote Java scripts for every reverse engineering step to perform decompiling automatically. The following main

<sup>2</sup><http://www.appchina.com>

<sup>3</sup><http://www.businessofapps.com/the-ultimate-app-store-list/>



Fig. 3. Reverse engineering for Android applications.

steps and resources which were used in decompiling the APK files to source code. 1)Unzip .apk files to get all .dex files. 2)Through scripts, we used Dex2jar<sup>4</sup> tool to get .jar files. 3)We used JD-CORE<sup>5</sup> decompiler to get .java source files from .jar files

### 4.3 Excluding Third-Party Libraries

Third-party libraries may affect the accuracy rate of malware detection and may also slow down the detection process. So, we excluded them during signature generation and detection process of malware. We uses a Whitelist to filter all third-party libraries. Although, it was not possible to build a complete Whitelist effectively because sometime obfuscation may change the name of packages. During experiments, we find many library files named "com/d/c/b", where come is the root folder. Which was quite hard to filter using Whitelist technique. So, We apply some filters against packages and classes names to ignore third-party library files. We extract the semantic features of different API calls and generated a vector sequences of directories. Then we perform comparison between these sequence vectors during malware detection and ignored these library files.

## 5. APK SIMILARITY (CODE CLONES) DETECTION

In this section, we explain how our proposed DroidSD clone detection tool works. Our proposed approach is a hybrid (semantic), which performs several screening filters to detect absolute and exact code clones in APK source code files. The clone detection method runs in a pipeline, where every process fully depends on the output of the previous step as shown in Fig. 4.

### 5.1 Preprocessing and Normalization

This is a first but very important process in APK code clone detection. In preprocessing code retrieved from source repository and split into tokens. We replace all integers, variables, functions, methods *etc.* into specified ids and numbers, meanwhile all import utilities, comments and spaces ignored as shown in first part of Fig. 5. Preprocessing of APK source code includes many steps, *i.e.* loading source code, verify and clean the source code, tokenization for lexical analysis.

### 5.2 Feature Extraction

This is very basic but the core phase of code clone detection in Android applications, main features from the source code of APK files would be extracted to transform source code into representation form for code comparison. These features include MD5 hash values of every chunk, path of source code file, first line number of every chunk and total lines of code in each file. After preprocessing and tokenization is done, the chunk formation performed on every file by putting 10 lines of code into a single chunk. The MD5 hashing used to get the hashing value of each chunk. To select some of these hashes

<sup>4</sup><https://github.com/pxb1988/dex2jar>

<sup>5</sup><https://github.com/nvniennot/jd-core-java>

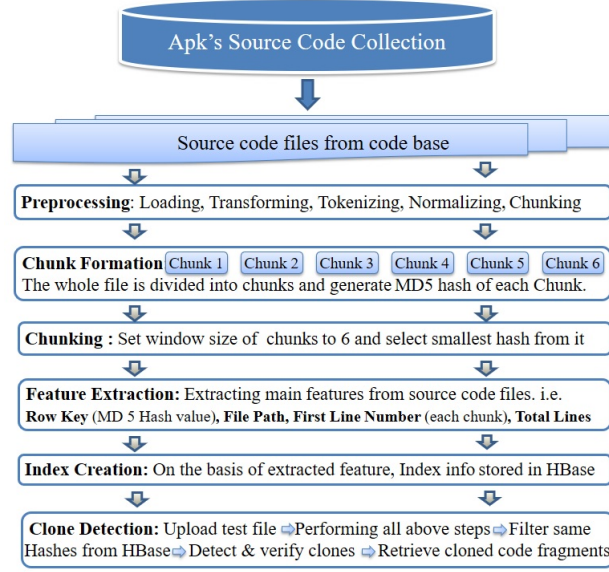


Fig. 4. DroidSD feature extraction and indexing process.

as a fingerprint from every source code file, we divide these hashes into groups (windows) and fixed the window size to 6 by choosing at least one fingerprint from each window as shown in Fig. 5 in red underlined hashes. By choosing at least one fingerprint from each window, we bound the maximum gap in fingerprints. Given a set of subject files, we want to find chunks matches between them should be satisfy two things.

1. If a chunk match, it verifies the actual threshold:  $t$ , then this chunk is detected.
2. There is no need to detect any chunk, who's threshold is smaller than threshold:  $u$ .

Where the  $t$  and  $u$  are constants and can be set by the user. The larger value of  $u$  makes us confident that the detected chunks are not coincidental.

Given a sequence of generated hash values  $h_1, h_2, \dots, h_n$ , if  $n > t - u$  then at least one hash  $h_i$  must be selected to guarantee the detection of all clone fragments of at least length  $t$ . By considering a simple approach, let the window size  $w = t - u + 1$ . Now consider a sequence of MD5 hashes  $h_1, h_2, \dots, h_n$ , which represents the chunks of a source code file. Every position  $1 \leq i \leq n - w + 1$  in this kind of sequence, defines a window of hash values  $h_i, \dots, h_{i+w-1}$ . So, it's mean to keep the guarantee of clone detection, it is very essential to select at least one hash value from every window size to be a fingerprint of the source code file. In our case we select the smallest hash value from every window as a fingerprint. The main reason behind selecting the minimum hash value is that the minimum hash value in one window is probably remain the minimum hash value in nearby windows, so the probability is that the minimum value of  $w$  (random numbers) is quite smaller than one supplementary random number. Therefore, many overlying windows select the same hash value and meanwhile the number of fingerprints selected values are very smaller than the total number of windows while still sustaining guarantee as shown in Fig. 5 in red underlined hashes.

The repeated hash values ignored by keeping it's information *i.e.* file name, starting line, ending line of that specific chunk. Each window size resulted in a term as a fingerprint, which later used for comparison of different source code files to detect code clones.

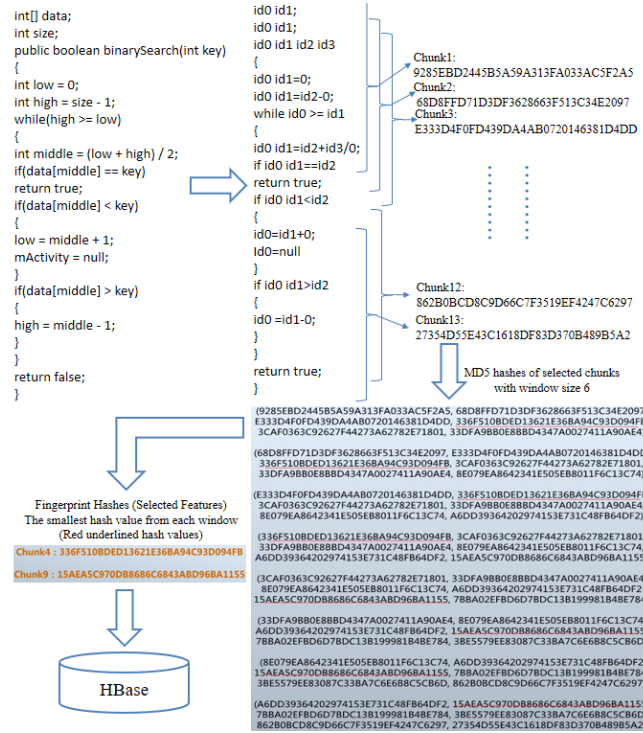


Fig. 5. Preprocessing, normalization, feature extraction.

Fingerprint value is directly proportional to code clones; it means that if the file has bigger fingerprint value, it has much affected by clones. The same method also used during code detection process, where hashes values of chunks later on compared with the index values stored chunks in HBase. Source code divided into chunks of equal sizes, each chunk consists of 10 lines of code, less than 15 LOC files discarded and ignored automatically. Chunk size is adjustable and can be varied any time, but for APK source code we set it at 10 LOC. In Chunk formation of a file, line from 1 to 10 would be formed as chunk1, line from 2 to 11 would be formed as chunk2, line from 3 to 12 would be formed as chunk 3 and so on. We add all these chunks into a chunk list and apply an MD5 hashing to get their hash values, which help to compare code fragments during clone detection. After getting all hashing values against all chunks, the least values from each window abstracted and saved in HBase as an index value. Here are the main extracted features.

**MD5 Hash:** This is the hashing value of each chunk.

**File Path:** Location where source code located in repository.

**First Line Number:** First line number of every chunk, to keep a record that from where this chunk starts. It helps us to retrieve clone fragments from the source code.

**Total Lines:** This feature is used to further verify that the cloned files have the same LOC or not. In case of both source files have the same number of LOC, it shows that file was fully copied from an original file.

### 5.3 Index Creation (Fingerprint Generation)

To minimize the indexing information and use less storage for index, there are two very important parameters considered in DroidSD: the fixed size of each hash value and

the size of each window. These two parameters determine the eigenvalue density and the accuracy of clone detection. In order to reduce the storage cost and to improve the detection accuracy, we have carried out several rounds of tuning the two parameters, and finally achieved satisfactory results. Since code cloning is typically done by rows, the first parameter is the size of the chunk, which computed by lines. Based on your own writing code, or by referring to the experience of open source code, the general 10 line can complete a simple functional unit, or a basically complete logical structure. So, we set the size of chunk to 10 lines. The second parameter is the size of the window (set of chunks) determines the shortest possible number of cloned lines of code, which set as 6.

All the extracted features from every source code file of APK saved in HBase as an index information. Our proposed index creation method is very fast and reliable. It just took 27 hours to build an index of 41 million files & 983 million LOC accurately. It considered takes very less time for large files as compare to small files of same size. Index function gets the source files path, total lines, chunks, MD5 hash values and stored in HBase. There were almost 83 million hash values were saved in HBase index table. The proposed index creation process is very fast, accurate, flexible and easy to maintainable.

#### 5.4 Clone Detection and Retrieval

In this section, we explained the final task of DroidSD, in which it filters the clone fragments from source code and perform the evaluation. The concept of MapReduce has been used to detect and retrieve the cloning files from repository against every subject system (to be detected). All the subject systems uploaded to the HDFS file system and become the input of map() in MapReduce. Then pass through certain criteria of preprocessing, normalization, feature extraction and chunk formation to detect clones in it as shown in Fig. 6.

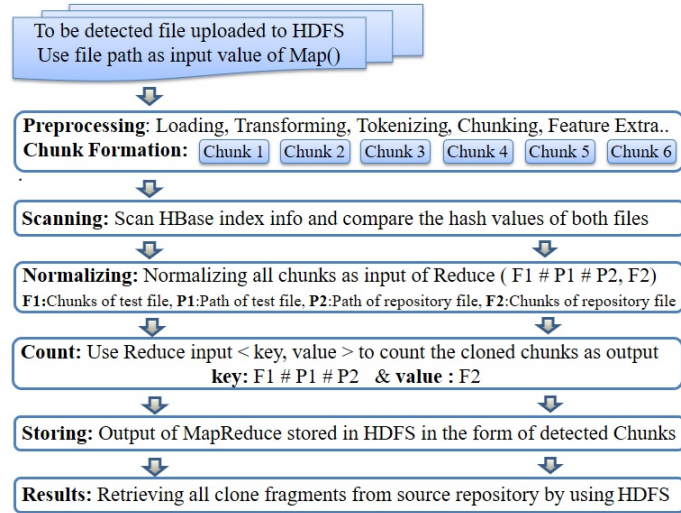


Fig. 6. Clone detection architecture.

At the initial stage of the detection process, all the preprocessing and normalization steps would have performed after uploading a subject system. After forming chunks of the subject system, the scanning process takes over the next procedure, which scans all the hash values of the index in HBase and compares them with the hash values of subject system. During the scanning process, it starts retrieving the chunks against similar



hash values from HBase repository. After scanning and retrieving all detected chunks, DroidSD performs normalization on all chunks through `reduce()`. Inputs of `reduce (F1# P1# P2, F2)` function are the values of test files, where F1 is chunks of the subject file, P1 is the path of the subject file, P2 is the path of repository files, and F2 is the chunks of repository files. Chunks of subject files and repository files further compared and evaluated by considering their path value. All detected chunks counted one by one to find their occurrence in different repository files. After performing count step, MapReduce retrieves all the chunks, which are detected as clones and stores into HDFS for further process of retrieving code segments.

### 5.5 Near-Miss (Type-3) Clones Detection

This is the distinguish characteristic of DroidSD to detect near-miss clones in Android application, where the source code have been changed by developers. These clone fragments have been changed by adding, modifying or removing statements during copying the source code from other applications. Many Type-3 clones have modifies by swapping statements in a source file or by combing multiple condition statements into one. Detecting such kind of similar code fragments from apk is very challenging for other Android clone detection approaches. So, we focused on detecting near-miss clones. Near-miss clones were further divided into different groups, *i.e.* very strong type, strong type, medium type and weak type and their detection results have been shown in Table 2. Fig. 9 shows one of the detected near-miss code fragment, which was detected in WhatsApp application. The other near-miss clone fragments have been discussed below in case study and the near-miss clones results have been shown in Table 4.

## 6. CASE STUDY OF CLONE DETECTION

In this part, the results of DroidSD approach have been shown and discussed. By downloading a big amount of Android applications from AppChina store as shown in Table 1, we transform all .apk files to .java source code files by using different reverse engineering techniques. We preprocessed and normalized almost 41 million files to build and save their index information in HBase. All Android application files decompiled in the same way and with the same tools, that's why the source code recovered by reverse engineering process was the same every time.

**Table 1. Source code and index info.**

Downloaded Apk	Decompiled Source files	Lines of Code	MD5 Hash Values
30,500	41,261,694	983,975,411	82,998,573

**System Specification:** Our APK clone detection approach was performed on Linux operating system. All process till clone retrieval was performed on a single machine (Intel core-i7, 3.60GHz\*8 & 24 GB of RAM) *i.e.* downloading Android apps, decompiling source code, preprocessing, normalizing, indexing and clone detection.

For the evaluation of DroidSD, we took top Android application named *Chrome*, *Firefox*, *Gmail*, *WhatsApp*, *GoogleMap*, *GooglePlayStore*, *Baidu* as a subject system to detect clones in them. These APK files were first decompiled through reverse engineering techniques to get their source code. HDFS has been used for uploading these applications decompiled source code to Hadoop file system for clone detection. Preprocessing, normalization, feature extraction were performed to make code ready for detection. Indexed

**Table 2. Detection Results of Chrome, Firefox, Gmail, WhatsApp, GoogleMap, Google-PlayStore, Baidu and BigCloneBench.**

App Name	Total Files	LOC	Detection Time	Type-1 Clones	Type-2 Clones	VST3 Clones	ST3 Clones	MT3 Clones	WT3,T4 Clones
Chrome	8,291	618,682	2 h 53 min	230,305	142,137	82,949	5,300	192,604	266,316
Firefox	3,614	578,124	2 h 38 min	286,850	212,861	294,305	7,938	543,236	610,174
Gmail	10,796	874,888	4 h 36 min	172,784	214,676	129,893	16,027	250,949	381,531
WhatsApp	3,895	628,426	3 h 10 min	484,323	169,346	277,656	190,956	512,977	256,093
GoogleMap	12,794	1,115,701	5 h 57 min	110,209	146,753	120,458	16,328	226,588	318,798
GooglePlayStore	11,133	899,689	5 h 3 min	205,095	234,390	129,430	13,409	219,788	385,042
Baidu	3,326	631,469	3 h 8 min	22,755	15,596	7,071	2,309	12,650	20,800
BigCloneBench in (IJaDataset)	51,499	10,431,956	7 h 15 min	91,387	1,572,672	451,319	103,827	759,867	1,693,033

```

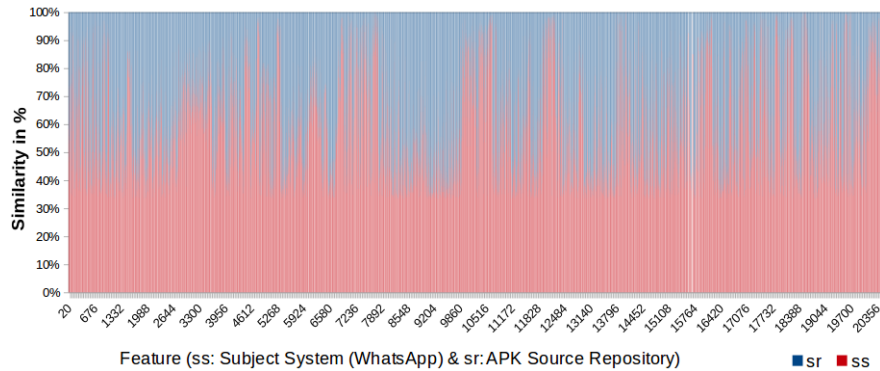
hadoop/BackUp/apk_collection/appChina/apk java source/Catwang1/com/google/android/gms/wearable/internal/zza.java 389
1903878 61shdfs://192.168.1.100:9000/sourcecode/12345678912345/whatsapp/android/support/v4/content/FileProvider.java$media/
hadoop/BackUp/apk_collection/appChina/apk java source/GoogleCalendar/android/support/v4/content/FileProvider.java 43
1903879 24shdfs://192.168.1.100:9000/sourcecode/12345678912345/whatsapp/com/google/android/gms/wearable/internal/a.java$media/
hadoop/BackUp/apk_collection/appChina/apk java source/ChromeBrowser/com/google/android/gms/t/xR.java 84
1903880 26shdfs://192.168.1.100:9000/sourcecode/12345678912345/whatsapp/com/google/android/gms/location/internal/e.java$media/
hadoop/BackUp/apk_collection/appChina/apk java source/Car-Net1/com/google/android/m4b/maps/r/e.java 20
1903881 2694shdfs://192.168.1.100:9000/sourcecode/12345678912345/whatsapp/com/google/android/m4b/maps/r/e.java 20

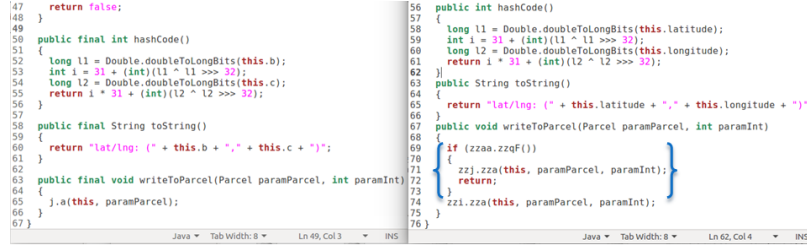
```

**Fig. 7. APK Clone Detection results (HDFS file).**

MD5 hash values of approximately 83 million from HBase were scanned and comparison between these hashes was performed through MapReduce. The detection time of each subject system is displayed in Table 2. Results of clone detection were retrieved from Reduce, which consists of feature similarity ratio in two files and their local path. Fig. 7 shows of only 3 selected results of subject system (*WhatsApp*), it is actually the HDFS file of size 417 MB, which consists of 1,903,894 retrieved clone results. Each result in it gives the actual path of detected cloned subject file (*WhatsApp*), repository files and the number of similar features in clone files. The 1903878 fully highlighted line in Fig. 7 is the file from WhatsApp, which was cloned in other Android applications in our source repository, names of detected clone Android applications in our source repository have been underlined with red color. The graphical representation of selected results of Fig. 7 has been shown in Fig. 8, which displays the actual similarity against different retrieved files. In Fig. 8, the files of subject system (*WhatsApp*) are in red color and the files of source repository are in blue color. The white lines shows that the features in these files are very less, *i.e.* the source code in these files is less than 10 LOC.

Table 2 shows the detection results of the subject systems against their total number of files, LOC, time and number of detected clones of Type-1, Type-2, Type-3. Types of clones actually depends on features similarity of files, *i.e.* similarity: 1.0 shows Type-1 clones, similarity: 0.9 - 1.0 shows Type-2 clones, similarity: 0.5 - 0.8 shows Type-3 clones. There is no hard and fast rule on when a clone is not syntactical similar, so it's

**Fig. 8. Features similarity of WhatsApp in APK source repository.**



Left picture: \$hdfs://192.168.../whatsapp/.../gms/maps/model/LatLng.java

Right picture: \$/.../appChina/apkSource/HellRider/unity/maps/ads/Unity.java

Fig. 9. Near-Miss clone detection result of APK.

**Table 3. Clones detected between Chrome and Firefox.**

App Name	Detection Time	Type-1 Clones	Type-2 Clones	VST3 Clones	ST3 Clones	MT3 Clones	WT3,T4 Clones
Chrome over Firefox	13 min	301	5085	670	221	1020	4820
Firefox over Chrome	11 min	67	2245	333	138	569	2181

quite hard to separate Type-3 and Type-4 clones, instead we divide and categorise them on the basis of their similarity measure, *i.e.* very strongly Type-3 (VST3) clones have the feature similarity of 80-90%, strongly Type-3 (ST3) clones have the feature similarity of 70-80%, moderately Type-3 (MT3) clones have the feature similarity of 60-70%, weakly Type-3/4 (WT3/4) clones have the feature similarity of 50-60%. The clones of less than 50% similarity can be Type-4 clones, but not surely evaluated and considered in this paper.

The result of Type-3 (near-miss) clone of our similarity detection approach has been shown in Fig. 9, which describe that new lines have been added to the source code of the original file, or may be the lines of source code were deleted from the original file. Furthermore, there were some function & method names have been changed in cloned file. Left picture in Fig. 9, is one of the subject file named *LatLng.java* from 3,895 files of *WhatsApp* application. The picture on the right is the file from another APK named *Hell-Rider*, which was detected and retrieved from our Android source code repository.

Another case study has been done for Batch clone detection (one to one), in which two systems have been used, *i.e.* Chrome and Firefox. The clone detection was automatically performed on each system one by one by considering one system as a subject system and another as a source repository and vice versa, the results are listed in Table 3. Through this experiment we can see the source code similarity between these two browser applications. It have been seen that many files of these application have been sharing the source code fragments. Through DroidSD we can even display those code fragments very effectively. Table 3 shows the results of Type-1 Type-2 and Type-3 clones Chrome over Firefox (Chrome as subject system) and Firefox over Chrome (Firefox as subject system).

## 7. EVALUATION

In this section, we have evaluated the performance of DroidSD. We have implemented our approach in real time environment to evaluate it's accuracy. To get the value of false negative and to make sure, if our APK code clone detection method is better enough, we supposed to count how many clones were not identified by DroidSD. So, we count the Recall and Precision. A perfect code similarity detection approach supposed to have Recall and Precision values both 100%. Recall is actually the fraction of all related files, which have been retrieved through a query. High Recall means that most of

the clones in that application have been found. In our approach, relevant clone files are those files which are supposed to be retrieved as cloned files, and retrieved files are those which were detected and retrieved as clones. In Precision, the relevant files retrieved by the query and it measures that how many irrelevant files were retrieved as a clone. High precision means that candidate clones are mostly actual and real clones. However, we could not able to find any benchmark for Android applications. Neither we could evaluate all the retrieved results. Because there were million clone results, which were retrieved from the source code of 30,500 million applications as shown in Table 2. So, we used an unbiased way to evaluate DroidSD through two ways. The first way by mutation framework and the second way BigCloneBench<sup>6</sup>.

**Mutation Framework:** The first way to evaluate DroidSD is creation a data-set by manual inserting different type of clone fragments from BigCloneBench in the source code files. We create our own data-set of 100 Android applications by inserting test code fragments of Type-1, Type-2 and Type-3 in these applications. We inserted 30 Type-1 clone fragments in first 30 applications, 30 Type-2 clone fragments in next 30 applications and 40 Type-3 clone fragments in last 40 applications. The data-set repository formed and ready for DroidSD evaluation. Then we extract features and build index of this data-set repository by using DroidSD. We take these inserted code fragments of BigCloneBench as a subject system and detect code clones in data-set repository. DroidSD detects clones at different threshold values to differentiate different types of clones. All retrieved clone fragments from data-set further manually evaluated by comparing them with the subject code fragments those inserted in data-set. Fig. 10 shows the clone detection results from data-set at different threshold values. As we minimize the threshold value, we able to detect Type-3 clones. The graph in Fig. 10 shows that DroidSD detects almost all inserted code fragments of Type-1 and Type-2 clones. But there is a small difference in detection of near-miss clones because of threshold value, if we minimize that value, we could detect almost all inserted fragments.

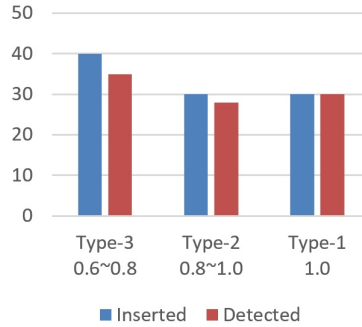


Fig. 10. Manually evaluated results of data-set.

**BigCloneBench:** BigCloneBench is a benchmark of manually collected and evaluated clone pairs in IJaDataset<sup>7</sup>. IJaDataset consists of almost 25,000 open source Java projects, 3 million files and 250 MLOC. BigCloneBench was built on the basis of IJaDataset and consists of almost 8 million validated clone pairs.

We executed DroidSD for IJaDataset and evaluated Recall with BigCloneBench (51,499 files, 10 MLOC). The results measured BigCloneBench are summarized per clone

<sup>6</sup><https://github.com/clonebench/BigCloneBench>

<sup>7</sup><https://jeffsvajlenko.weebly.com/bigcloneeval.html>

```

2983995 29$hdfs://192.168.1/bcb_reduced/2/selected/1005719.java$media/dataSet/IJaDataSet_0_9/files/1005166.java 45
2983996 29$hdfs://192.168.1/bcb_reduced/2/selected/1005719.java$media/dataSet/IJaDataSet_0_9/files/1006102.java 150
2983997 29$hdfs://192.168.1/bcb_reduced/2/selected/1005719.java$media/dataSet/IJaDataSet_0_9/files/1008122.java 57
2983998 29$hdfs://192.168.1/bcb_reduced/2/selected/1005719.java$media/dataSet/IJaDataSet_0_9/files/101305.java 45
2983999 29$hdfs://192.168.1/bcb_reduced/2/selected/1005719.java$media/dataSet/IJaDataSet_0_9/files/1014216.java 45
2984000 29$hdfs://192.168.1/bcb_reduced/2/selected/1005719.java$media/dataSet/IJaDataSet_0_9/files/1018489.java 63
2984001 29$hdfs://192.168.1/bcb_reduced/2/selected/1005719.java$media/dataSet/IJaDataSet_0_9/files/1019260.java 45

```

Fig. 11. BigCloneBench results from IJaDataSet (HDFS file).

type in Table 2 and the HDFS result file have been shown in Fig. 11. DroidSD has perfectly detected Type-1 and Type-2 clones in BigCloneBench. The VST3 clones are also have excellent detection rate as shown in Table 4. To detect near-miss code fragments of Type-3 was not possible in previous Android similarity detection techniques. The weak Type-3 clones which can be Type-4 clones were not evaluated in this paper, So we consider it as a 0 Recall and Precision.

As for Recall, there exist a high quality benchmark, but to count Precision, it was still an open problem. So, we used manually evaluation method to count Precision for DroidSD. We randomly selected 200 files from the results and for fair evaluation, we divide them on three judges, who were having the knowledge of source code similarity measures. After combining their results, we found that DroidSD has high Precision (97%) for Type-1 clones, (91%) for Type-2, (83%) for very strong Type-3 clones.

**Table 4. Recall and precision (BigCloneBench).**

Clone Types	Type-1	Type-2	VST3	ST3	MT3	WT3
Recall (%)	98	93	87	63	30	0
Precision (%)	97	91	83	49	17	0

During evaluation process, we have performed full application similarity detection and partial application similarity detection. Using full application similarity detection, we detected all clone files from dataset repository against our all subject systems as the results are shown in Table 2. It has been detected that two same applications with different versions have much more common features than different application. In the case of partial similarity detection, our approach has successfully found the applications which share some source files or part of their source code. our approach also retrieves and count all features in the subject file and the clone file, if the total number of features in both files are same, it's also the symbol of 100% similarity in files. Both full and partial similarity detection require finding similar code fragments in application's source code.

For further evaluation and to check the accuracy, we perform manual analysis of clone fragments. We manually evaluated almost 200 code pairs from Chrome (subject system), which were reported to be cloned. To manually evaluate every clone fragment, we opened both files side by side and verify if they were similar. Manually evaluation on clone fragments performed for Type-1, Type-2 and Type-3 code clones. Meanwhile, Recall and Precision counted against every clone type as shown in Table 5.

**Table 5. Recall and precision measurement (Chrome apk).**

Clone Types	Type-1	Type-2	VST3	ST3	MT3
Evaluation Files	20	20	20	20	20
Recall (%)	93	86	75	54	23
Evaluation Files	20	20	20	20	20
Precision (%)	96	89	81	32	17

**Scalability:** The scalability of DroidSD was evaluated by using different inputs with varying the size of LOC, at different level of granularities. In default execution time scales with the size of input (LOC). As our approach is based on Hadoop, HBase and MapReduce, so the scalability is not an issue. DroidSD can handle million of files and billion of LOC very effectively. IJaDataset and BigCloneBench are considered two big data-sets, which we used for the evaluation of DroidSD. For much large scale systems, our technique can be extended to multiple clusters which can handle any size of input. The results in Tables 2-5 proves the performance ability of our approach.

## 8. COMPARISON WITH EXISTING APPROACHES

Most of the application's clone and similarity detection approaches are simple hashing, semantic feature based, PDG based, API method based and UI based and even they were unable to detect near-miss clone fragments from apk source code. We considered a few detection approaches to compare with DroidSD, *i.e.* PDG based approaches are not scalable and can not be implemented on clone detection on a market with over million applications. Kim's technique (API method) executes applications and then collect all API call sequences as birthmarks. API method can not work when apps are encrypted by Ijimi because API method would not be able to get exact API traces of the application. Soh (UI method) uses the user interface information to check the similarity. If some fake activities or code fragments inserted in apps, which does not influence or make any change in user interface, UI method can not detect these kind of similar code fragments. WuKong [26] detects code clone by using feature matrices for each app by building  $n*m$  Characteristic Matrix, which is an abstraction of code segments. DroidSD detects similarity of Android applications at source code level by extracting it's main features.

## 9. LIMITATION

DroidSD can detect code clones in Android applications very effectively but it's still hard to differentiate between the clone code and the original code. DroidSD can detect Type-1,2,3 code clones with high accuracy rate but it's very hard to detect Type-4 clones, because it's a token based clone detection approach and up to our knowledge there is no any token based approach which can detect Type-4 clones but we considered it as a probability that the clones of less than 50% similarity can be Type-4 clones, but not surely evaluated and considered in this paper.

## 10. CONCLUSION AND FUTURE WORK

In this paper, we have proposed DroidSD a novel approach for clone detection in Android applications at source code level. Reverse engineering has been used for de-compiling the source code of Android applications. Several tasks have been performed to preprocess the source code *i.e.* cleaning, transforming, normalizing and tokenizing *etc.* Index of each source code file was built by extracting their main selected features. In total, index of 982 million LOC source code was built just less than 3 days. DroidSD can detect Type-1, Type-2 and Type-3 code clones from any kind of Android application, which helps us to check the fake, vulnerable and malware application. Our clone detection approach is incrementable and can be extended to distributed clone detection approach,

where we can build an index of a large amount of Android applications and detect code clones from apps on large scale. DroidSD can be used by Android app stores to check the accuracy, validity and repackaging similarity in every application before upload into app stores. For future concerns, we are planning to develop an algorithm, which can help us to detect Type-4 clones effectively, and meanwhile detect the malicious or vulnerable code segments from Android applications to overcome the security threat.

## ACKNOWLEDGEMENT

This research done in Key Laboratory of Information System Security, School of Software, Tsinghua University Beijing, China and it was supported by the National Natural Science foundation of China under Grant No: 61540020.

## REFERENCES

1. K. Allix, T. F. Bissyandé, J. Klein, and Y. L. Traon, "Androzoo: Collecting millions of android apps for the research community," in *Proceedings of IEEE/ACM 13th Working Conference on Mining Software Repositories*, 2016, pp. 468-471.
2. J. Chen, M. H. Alalfi, T. R. Dean, and Y. Zou, "Detecting android malware using clone detection," *Journal of Computer Science and Technology*, Vol. 30, 2015, pp. 942-956.
3. C. Amrutkar, P. Traynor, and P. C. Van Oorschot, "An empirical evaluation of security indicators in mobile web browsers," *IEEE Transactions on Mobile Computing*, Vol. 14, 2015, pp. 889-903.
4. J. Akram, Q. Liang, and L. Ping, "Vcpr: Vulnerable code is identifiable when a patch is released (hacker's perspective)," in *Proceedings of the 12th IEEE Conference on Software Testing, Validation and Verification*, 2019, pp. 402-413.
5. J. Akram and P. Luo, "How to build a vulnerability benchmark to overcome cyber security attacks," in *IET Information Security*, 2019,
6. J. Crussell, C. Gibler, and H. Chen, "Andarwin: Scalable detection of android application clones based on semantics," *IEEE Transactions on Mobile Computing*, Vol. 14, 2015, pp. 2007-2019.
7. Y. Y. Ng, H. Zhou, Z. Ji, H. Luo, and Y. Dong, "Which android app store can be trusted in china?" in *Proceedings of IEEE 38th Annual Computer Software and Applications Conference*, 2014, pp. 509-518.
8. E. Juergens, F. Deissenboeck, B. Hummel, and S. Wagner, "Do code clones matter?" in *Software Engineering, 2009. ICSE 2009. IEEE 31st International Conference on*. IEEE, 2009, pp. 485-495.
9. K. Chen, P. Liu, and Y. Zhang, "Achieving accuracy and scalability simultaneously in detecting application clones on android markets," in *Proceedings of the 36th ACM International Conference on Software Engineering*, 2014, pp. 175-186.
10. H. Niu, T. Yang, and S. Niu, "Clone analysis and detection in android applications," in *Proceedings of IEEE 3rd International Conference on Systems and Informatics*, 2016, pp. 520-525.
11. J. Akram, Z. Shi, M. Mumtaz, and P. Luo, "Droidcc: A scalable clone detection approach for android applications to detect similarity at source code level," in *Proceedings of IEEE 42nd Annual Computer Software and Applications Conference*, 2018, pp. 100-105.

12. J. Akram, M. Mumtaz, J. Gul, and P. Luo, "Droidmd: An efficient and scalable android malware detection approach at source code level," *International Journal of Information and Computer Security*, Vol. 11, 2019.
13. S. Alam, R. Riley, I. Sogukpinar, and N. Carkaci, "Droidclone: Detecting android malware variants by exposing code clones," in *Proceedings of IEEE 6th International Conference on Digital Information and Communication Technology and its Applications*, 2016, pp. 79-84.
14. A. Sheneamer and J. Kalita, "A survey of software clone detection techniques," *International Journal of Computer Applications*, Vol. 137, 2016, pp. 0975-8887.
15. Y. Dang, D. Zhang, S. Ge, R. Huang, C. Chu, and T. Xie, "Transferring code-clone detection and analysis to practice," in *Proceedings of IEEE/ACM 39th International Conference on Software Engineering: Software Engineering in Practice Track*, 2017, pp. 53-62.
16. S. Schleimer, D. S. Wilkerson, and A. Aiken, "Winnowing: local algorithms for document fingerprinting," in *Proceedings of ACM SIGMOD International Conference on Management of Data*, 2003, pp. 76-85.
17. M. A. Nishi and K. Damevski, "Scalable code clone detection and search based on adaptive prefix filtering," *Journal of Systems and Software*, Vol. 137, 2018, pp. 130-142.
18. A.-F. Mubarak-Ali, S. Sulaiman, S. M. Syed-Mohamad, and Z. Xing, "Code clone detection and analysis in open source applications," *Computer Systems and Software Engineering: Concepts, Methodologies, Tools, and Applications*, 2018, pp. 1112-1127.
19. A. Ghosh and Y. Lee, "An empirical study of a hybrid code clone detection approach on java byte code," *GSTF Journal on Computing*, Vol. 5, 2017, pp. 34-45.
20. A. Gupta and B. Suri, "A survey on code clone, its behavior and applications," *Networking Communication and Data Knowledge Engineering*, 2018, pp. 27-39.
21. J. Akram, Z. Shi, M. Mumtaz, and P. Luo, "Dccd: An efficient and scalable distributed code clone detection technique for big code," in *Proceedings of the 30th International Conference on Software Engineering and Knowledge Engineering*, 2018, pp. 354-359.
22. X. Cheng, H. Zhong, Y. Chen, Z. Hu, and J. Zhao, "Rule-directed code clone synchronization," in *Proceedings of IEEE 24th International Conference on Program Comprehension*, 2016, pp. 1-10.
23. W. Zhou, Y. Zhou, X. Jiang, and P. Ning, "Detecting repackaged smartphone applications in third-party android marketplaces," in *Proceedings of ACM Conference on Data and Application Security and Privacy*, 2012, pp. 317-326.
24. S. Kim, E. Kim, and J. Choi, "A method for detecting illegally copied apk files on the network," in *Proceedings of ACM Research in Applied Computation Symposium*, 2012, pp. 253-256.
25. J. Svajlenko and C. K. Roy, "Cloneworks: a fast and flexible large-scale near-miss clone detection tool," in *Proceedings of the 39th IEEE International Conference on Software Engineering Companion*, 2017, pp. 177-179.
26. H. Wang, Y. Guo, Z. Ma, and X. Chen, "Wukong: A scalable and accurate two-phase approach to android app clone detection," in *Proceedings of ACM International Symposium on Software Testing and Analysis*, 2015, pp. 71-82.





**Junaid Akram** is a member of IEEE Computer Society. He received his 1st Master degree in major of Information Technology from Pakistan in year 2009, and 2nd Mater degree in major of Communication and Information System from China in year 2015. Recently, he is Ph.D. scholar in School of Software at Tsinghua University China. His current research interests include software reuse, information security, vulnerability detection and code clone detection.



**Zhendong Shi** got his B.S. degree from Yunan University China. Now he is studding a master degree program at Tsinghua University China of major entitled Software Engineering. His research interests include big data, software verification, vulnerability analysis and parallel computing.



**Majid Mumtaz** is a faculty member of COMSATS Institute of Information Technology Pakistan. He is doing Ph.D. in Tsinghua University Beijing China. His primary research interests are Information and communication systems security, cryptography, algebraic cryptanalysis and mobile application security.



**Luo Ping** is a faculty member in Tsinghua University, China. He received Ph.D. degree in Applied Mathematics from Chinese Academy of Sciences Institute of Systems Science in year in 1996. He joined Department of Computer Science and Technology, Tsinghua University in 2008. His research interests include information security, cryptography, vulnerability analysis and attack, database vulnerabilities and security.