# LACTA: An Enhanced Automatic Software Categorization on the Native Code of Android Applications

Cheng-Zen Yang and Ming-Hsuan Tu

*Abstract*—**Since Android has become a popular software platform for mobile devices recently, many software repositories collect and archive Android applications to facilitate the dissemination of Android applications. As the number of new Android applications tends to be rapidly increased in the near future, automatic software categorization will be in great demand. Although there are many approaches proposed for automatic software categorization, they do not consider the challenges specifically for Android applications. In this paper, we propose an enhancement called LACTA based on LACT to tackle this problem. LACTA extensively employs Android domain knowledge in the process and uses LDA to extract meaningful software topics for classification. We have conducted empirical experiments with 42 applications. The experimental results show that LACTA has promising improvements under the consideration of Android domain knowledge.**

*Index Terms*—**automatic software categorization, latent Dirichlet allocation, Android, information retrieval**

## I. INTRODUCTION

AS Android has become a popular software platform for mobile devices recently, many software repositories for Android applications like Android Market (market.android.com) are also populated accordingly. For ease of browsing and searching for related software, they categorize the archived Android applications into groups as traditional software repositories for shareware and open-source software, such as SourceForge.net, However, the classification is usually determined manually by users or administrators for most Android repositories. As the number of new Android applications tends to be rapidly increased in the near future, automatic software categorization will be in great demand for management of Android application archives.

The research on automatic software categorization has been discussed for many years. However, most past studies mainly focus on the problem of classifying software components to facilitate software reuse [1]-[3]. In addition,

Cheng-Zen Yang is with the Department of Computer Science and Engineering, Yuan Ze University, Taoyuan, Taiwan. (phone:+886-3-4638800 ext 2361; e-mail: czyang@saturn.yzu.edu.tw)
Ming-Hsuan Tu is with the Department of Computer Science and Engineering, Yuan Ze University, Taoyuan, 32003 Taiwan. (e-mail: s973350@mail.yzu.edu.tw)

many proposed approaches for automatic software categorization need text information from associated documents or comments in the source code [1]-[5]. Although recently several approaches, such as MUDABlue [6], [7] and LACT [8], have been proposed to address the issue of automatically classifying software applications fully based on their source code, there are still challenges in applying these approaches to Android applications.

First, most disseminated Android applications are distributed in the native Android application package (APK) from. Therefore, the APK files need to be first decompiled to obtain the Java source code for further classification. However, many meaningful identifiers in the original Java source code will be converted to be meaningless in the decompilation process. Therefore, the mount of available useful textual information is much less in the decompiled Android source code than in other open-source software projects. To the best of our knowledge, no previous work has considered this difficulty in the classification design.

Second, there are many XML files describing accompanied resources in Android applications. How to employ these XML files to enhance the classification performance needs to be specifically considered for Android applications.

In this paper, we propose an enhancement based on the LACT approach [8] to improve the automatic software categorization performance for Android applications. The enhancement is thus called LACTA (LACT for Android) in which Latent Dirichlet Allocation (LDA) [9] is employed to extract meaningful software topics by extensively using Android domain knowledge. Then the software categories can be automatically determined from the extract topics. Thereafter the applications can be classified into the determined categories based on their identified topics.

To evaluate the effectiveness of LACTA, we collected 42 Android applications from the Internet and conducted empirical experiments in which LACTA was compared with the original LACT. With the enhancement of Android domain knowledge, LACTA outperforms LACT in classification performance. In addition, LACTA also extract more meaningful topic words for category generation.

The rest of the paper is organized as follows. In Section2, we briefly review precious research work on automatic software categorization. Section 3 describes the design details of LACTA. Section 4 presents the experiments and provides discussion for the experimental results. In Section 5, we discuss the potential threats to the validity of experimental results. Finally, Section 6 concludes the paper and describes

our future wok.

## II. RELATED WORK

Automatic software categorization obtains notifications when the information retrieval (IR) techniques are applied to the software engineering domain. In 1991, a categorization tool called GURU was proposed to automatically classify software libraries by analyzing the associated documents and manuals with IR techniques [1]. In GURU, documents associated with software libraries, such as manual pages, are indexed and classified into clusters using basic IR techniques like removing less significant features and similarity calculation. However, the classification of GURU relies on the manually determined categorical hierarchy. In addition, GURU does not consider the program code in its IR processing.

In 1995, Merkl addressed the issue of organizing software components using Self-Organizing Neural Networks [2]. In this approach, both the ART model and the self-organizing map (SOM) are used to explore the semantic similarities. As GURU, this approach extracts keywords of software components from manuals rather than the program code. Therefore, this approach has the same shortcomings as GURU. In addition, it only considers the classification at the component level, rather than the software level. A similar approach also using SOM is proposed in [3] by Chan and Spracklen. Without the assistance of manual information, this approach analyzes the source code for the classification. However, it still only considers the component level classification.

In 2002, a supervised machine learning approached was proposed to use support vector machine (SVM) classifiers for automatic software categorization [4], [5]. With the prominent classification performance of SVM on text documents, this approach achieves an average accuracy of 43% in software categorization. However, in this approach text information of comments and README documents plays an important role. In addition, the categories are fixed due to the characteristics of supervised learning.

In 2004, MUDABlue was proposed to tackle the software categorization problem using an unsupervised approach based on Latent Semantic Analysis (LSA) [10] technique [6], [7]. With the effectiveness of LSA in extracting latent semantics from source code, MUDABlue has three major advantages. First, the classification categories of software projects can be automatically constructed. Second, MUDABlue analyzes source code only. It does not require additional information of software manuals or comments. Third, a software project may be classified into several categories rather than a single category as in the traditional classification work. Using 41 C programs as the dateset, MUDABlue shows its promising performance [6], [7]. However, a later study shows that some category names are difficult to interpret in MUDABlue [8].

In 2007, Kuhn, Ducasse, and Girba proposed a technique called *Software Clustering* (SC) which also uses LSA to find software topics [11]. In SC, source code and comments are analyzed to extract linguistic information. However, it is mainly designed for software having the original source code.
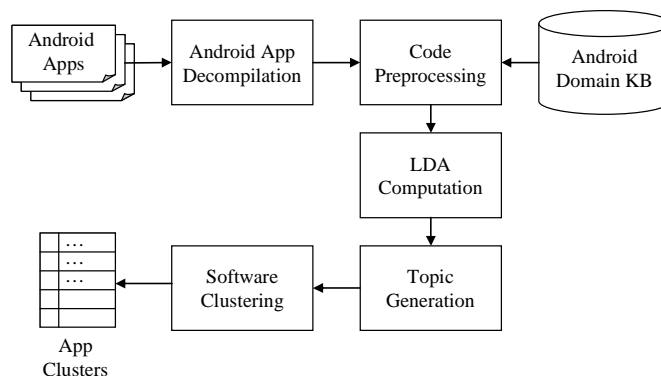


Fig. 1. The processing flow in LACTA for automatic software categorization on Android applications.

In addition, it focuses on clustering software applications according to the vocabulary, rather than classifying them according to the category meanings.

In 2009, Tian, Revelle, and Poshyvanyk proposed LACT [8] for software categorization using Latent Dirichlet Allocation (LDA) which is a more advanced topic extraction technique [9]. Although LACT can only achieve comparable performance as MUDABlue, it can handle the divergences of different programming languages. However, the lack of extensively employing domain knowledge of software platforms, its performance is limited for Android applications as shown in our experiments.

## III. LACTA DESIGN

For automatic software categorization on Android applications, we enhance LACT by extensively including Android platform domain knowledge in the classification process. Fig. 1 illustrates the processing flow in LACTA. The details of each step is elaborated in the rest of this section.

### A. Android Application Decompilation

The objective of LACTA is to directly handle native code of Android applications. Since Android applications are packed into Android application package (APK) files for dissemination, LACTA performs decompilation to convert these APK files into Java source code for the following

```
package ru.gelin.android.weather.notification.skin;
...
import android.app.NotificationManager;
import android.app.PendingIntent;
...
public abstract class WeatherNotificationReceiver extends BroadcastReceiver
{
 public static final String ACTION_WEATHER_UPDATE;
 public static final String EXTRA_ENABLE_NOTIFICATION;
 public static final String EXTRA_WEATHER;
...
 protected static PendingIntent getMainActivityPendingIntent(Context paramContext)
  {
   Intent localIntent1 = new Intent(paramContext, MainActivity.class);
   Intent localIntent2 = localIntent1.addFlags(536870912);
   return PendingIntent.getActivity(paramContext, 0, localIntent1, 0);
  }
...
   }
```

Fig. 2. A decompilation results for an application weatherNotification 0.1.2.

processing. For decompilation, LACTA first uses the *dex2jar* tool [12] to convert APK files into jar files, and then uses the *JD-GUI* Java decompiler [13] to obtain the Java source code. Fig. 2 is a decompilation example in which we get the Java source code for an application weatherNotification 0.1.2. In this step, LACTA also uses android-apktool [14] to get the XML resource files in the Android applications. From the figure, we can also find that many local variables, such as *localIntent1*, lose their meanings after the decompilation. Only identifiers declared as "public", "private", and "protected" or declared as class names have their original meanings.

### B. Source Code Preprocessing

In this step, traditional information retrieval (IR) preprocessing techniques are applied to the decompiled Java source code. A large amount of Android domain knowledge is also applied in the preprocessing step to obtain high-quality features for classification. First, the non-literal characters and the reserved keywords are removed from the source code, because these terms have little meaning for software categorization.

Then the identifiers in the source code are split to extract terms for the classification work. For example, "gameStart()" is split into "game" and "start", and "music_player" is split into "music" and "player". The splitting rules used in LACTA are similar to the rules in [15]. The characters like underscore and hyphen are used as the delimiters.

There are two main reasons to perform split operations on identifiers. First, the decompiled Java code does not contain comment information. Therefore, the split operations can extract more meaningful terms from identifiers. Second, many identifiers may contain common terms such as "get" in "getActiveNetworkInfo()". These common words appear in many applications and need to be eliminated to improve the discriminability of the following classification work. However, elimination of common terms influences the classification performance because preserving some common terms may be contributive to the classification accuracy. For example, the words "player" and "video" are representative for multimedia applications. Therefore, a term will be removed only if it appears in more than one-half of applications, and does not appear in the local string resource files, such as Strings.xml.

After term extraction, common IR techniques, such as stopword removal and stemming, are applied to these words. Then each Android application is represented by a collection of preprocessed terms for LDA computation. In addition, the numbers of the terms appearing in both Java code and the string resource files are multiplied by a weighting factor $\alpha$ because these terms may have significant categorical meaning. This weighting scheme is equivalent to $\alpha \times tf_i$, where $tf_i$ is the term frequency of term $i$.

### C. LDA Computation

Latent Dirichlet Allocation (LDA) has been proven to be an effective mechanism to mine the latent semantic topics from documents [9]. In this step, we follow the design of LACT [8] to extract the word-topic distribution matrices and the topic-software distribution matrices using LDA. Therefore,

each application can have multiple classification topics each of which consist of a collection of terms. These two kinds of matrices are used for further topic generation and software clustering.

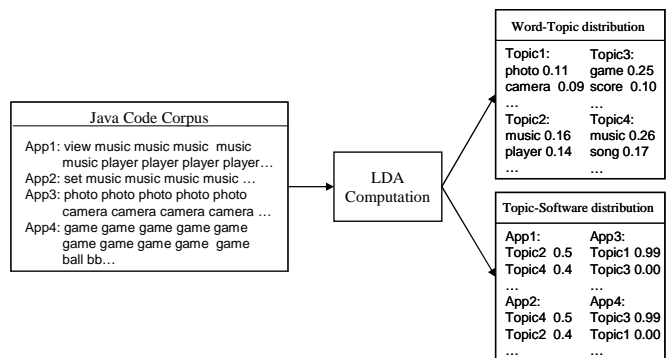Since the number of topics is an adjustable parameter in



Fig. 3. The LDA computations generate the word-topic distribution matrices and the topic-software distribution matrices.

LDA, different settings are evaluated in the experiments. Fig. 3 shows an example in which there are four applications and each application has a topic rank list for four topics extracted in LDA computations.

### D. Topic Generation

To automatically decide the classification categories, topics of similar semantics will be clustered into a category from the word-topic matrices. In this step, topic $t_1$ and $t_2$ are said to be semantically similar if the cosine similarity of them is large than a predefined threshold $h_t$. The cosine similarity is computed as follows:

$$similarity\ (t_1, t_2) = \frac{t_1 \bullet t_2}{|t_1| \times |t_2|} \qquad (1)$$

where topic $t_1$ is represented by the vector of extracted topic terms in the word-topic matrices and so is $t_2$. If the cosine similarity is large than $h_t$, these two topics are merged. Another topic $t_i$ will be merged into this cluster only if $t_i$ is similar to all existing topics in the cluster. Then for each category there is a corresponding topic cluster containing one or more topics.

Although this simple clustering process is effective in most cases, a *topic drifting* problem may exist when a topic cluster has two topics that are not similar, i.e., one topic is drifted far away from another topic. To mitigate the topic drifting problem in our approach, we use $h_t = 0.9$ in topic generation.

### E. Software Clustering

In the final step, software applications are classified into the topic clusters according to a predefined topic-software threshold $h_c$, which defines how many topic clusters will be considered for an application. If an application $App_j$ has a topic $t_m$ whose distribution value is large than $h_c$, it will be classified into the corresponding category.

## IV. CLASSIFICATION EVALUATION

To study the classification performance of LACTA, we collected 42 Android applications from the Internet. Since we collected the applications from several sites and each site has

TABLE I
THE COLLECTED 42 ANDROID APPLICATIONS FOR CLASSIFICATION
STUDY.

| Category | Applications |
|---|---|
| Battery | BatteryBar, BatteryChecker, One_Touch_Battery_Saver |
| Camera | Camera_360_0, Camera_Fun_Free_2, Camera_Illusion_1 |
| Communication | PhoneQ_Lite_1, ReChat_0, Twitter_2 |
| Finance | AndTip_1, Auto_Loan_Calculator_1 |
| Game | Air_Control_Lite_1, AirAttack_Lite_3, Andoku_1, Balance_The_Beer_1, BallDroppings_Lite_1, Basketball_Shot_1, BasketBall_v1, City_Jump_1, Jumper0, Paddle_Bounce_1, Pro_Basketball_Scores_2, Sand_Blaster_1, Toss_It_1 |
| Multimedia | Adobe_Flash_Player_10, AMPlayer_0, DAAP_Media_Player_0, Google_music, Movies_2, Music_Queue_1, MusicCube_1, SPB_TV_lite_1, Tranquilize_v1, TV_Listings_2, Video_Player_1, YouTube_2, |
| Reader | Adobe_reader_2, BeamReader |
| Weather | Animated_Weather_Free_2, CityWeather_1, Weather_notification_0, Windfinder_1 |

TABLE II
THE CLASSIFICATION PERFORMANCE OF PRECISION AND RECALL FOR
LACTA.

| Topics | Threshold $h_c$ | | | | |
|---|---|---|---|---|---|
| | 0.01 | 0.02 | 0.05 | 0.1 | 0.2 |
| 10 | 0.38, 0.62 | 0.42, 0.62 | 0.46, 0.60 | 0.47, 0.60 | 0.48, 0.57 |
| 20 | 0.66, 0.83 | 0.67, 0.81 | 0.72, 0.79 | 0.74, 0.76 | 0.74, 0.74 |
| 30 | 0.78, 0.86 | 0.77, 0.83 | 0.79, 0.83 | 0.78, 0.83 | 0.77, 0.81 |
| 40 | 0.84, 0.98 | 0.84, 0.98 | 0.85, 0.98 | **0.89, 0.98** | 0.91, 0.95 |
| 50 | 0.87, 0.93 | 0.87, 0.93 | 0.87, 0.93 | 0.87, 0.93 | 0.86, 0.88 |
| 60 | 0.87, 0.93 | 0.89, 0.93 | 0.89, 0.93 | 0.89, 0.93 | 0.91, 0.93 |
| 70 | 0.85, 0.90 | 0.87, 0.90 | 0.86, 0.90 | 0.86, 0.90 | 0.86, 0.90 |
| 80 | 0.75, 0.95 | 0.78, 0.95 | 0.80, 0.95 | 0.87, 0.95 | 0.91, 0.95 |

TABLE III
THE CLASSIFICATION PERFORMANCE OF PRECISION AND RECALL FOR LACT.

| Topics | Threshold $h_c$ | | | | |
|---|---|---|---|---|---|
| | 0.01 | 0.02 | 0.05 | 0.1 | 0.2 |
| 10 | 0.15, 0.40 | 0.17, 0.40 | 0.17, 0.31 | 0.18, 0.29 | 0.19, 0.24 |
| 20 | 0.17, 0.45 | 0.20, 0.40 | 0.23, 0.36 | 0.29, 0.33 | 0.29, 0.33 |
| 30 | 0.22, 0.52 | 0.22, 0.52 | 0.22, 0.52 | 0.22, 0.50 | 0.23, 0.45 |
| 40 | 0.19, 0.57 | 0.20, 0.57 | 0.19, 0.52 | 0.21, 0.52 | 0.24, 0.50 |
| 50 | 0.18, 0.64 | 0.19, 0.64 | 0.21, 0.64 | 0.23, 0.64 | **0.27, 0.57** |
| 60 | 0.20, 0.48 | 0.19, 0.48 | 0.21, 0.48 | 0.21, 0.48 | 0.20, 0.40 |
| 70 | 0.19, 0.48 | 0.19, 0.48 | 0.20, 0.48 | 0.20, 0.48 | 0.20, 0.40 |
| 80 | 0.21, 0.55 | 0.21, 0.55 | 0.22, 0.55 | 0.23, 0.55 | 0.24, 0.48 |

its own classification, we manually classified the 42 applications into 8 categories. Table 1 describes the information of the dataset used in the experiments.

In the experiment, we used GibbLDA++ [16] as the LDA computation engine. The weighting factor $\alpha$ was 5, the threshold $h_t$ was 0.9, and the topic-software threshold $h_c$ ranged from 0.01 to 0.2 to discuss the impact of different topic-software thresholds.

To evaluate the effectiveness of LACTA, we also implemented LACT as the baseline. In the experiments, *precision* and *recall* was computed as the performance metrics for both LACT and LACTA. The classification correctness is manually decided by inspecting the topic words and the application contents as in the original LACT work [8]. Therefore, the generated categories are manually mapped to the categories of Table 1 for performance evaluation. The precision measure is the fraction of the number of applications correctly categorized divided by the number of all applications categorized. The recall measure is the fraction of the number of applications correctly categorized divided by the number of applications belonging to that category. In the table, we also mark the results with the highest $F_1$ measure, where $F_1 = (2 \times precision \times recall)/(precision + recall)$.

Table 2 shows the precision and recall results for LACTA. In the table, we can find that the topic-software distribution threshold $h_c$ and the number of topics influence the performance significantly. When the number of topics is 40 and $h_c$=0.1, LACTA has the best precision and recall.

Table 3 shows the precision and recall results for LACT. The experimental results show that LACT has the best $F_1$ performance when the topic number is 50 and $h_c$=0.2. However, LACT cannot outperform LACTA in all cases. There are two main reasons. First, the original LACT does not consider the domain knowledge of Android programming. Therefore, many platform-related topics will be generated in LACT. For example, "Activity" and "Resource" are two common words in Android programming, and their

appearances lower the classification performance. Second, the original design of LACT focuses on open-source software applications whose source code is accessible. However, only decompiled source code is available for Android applications. Therefore, the significance of many ambiguous identifiers such as "a" and "aa" in the decompiled Java code needs to be adjusted. In LACTA, we extensively utilize the string resource files in Android applications to adjust the weights of the extracted terms.

Using the settings of the highest $F_1$ measure, LACTA automatically decides 33 categories for these 42 applications. Compared to the 8 categories defined in Table 1, LACTA can find more categories with more specific meanings. Table 4

TABLE IV
SOME AUTOMATICALLY GENERATED SOFTWARE CATEGORIES IN LACTA.

| Category ID | Topic Words | Applications |
|---|---|---|
| 1 | file, playlist, prefer, artist, album, audioformat | AMPlayer_0, Adobe_Flash_Player_10 |
| 2 | battery, widget, long, state | BatteryChecker, One_Touch_Battery_Saver |
| 3 | song, player, music, listen, headset | Music_Queue_1, musicCube_1 |
| 4 | plusmo, widget, drop, ball, new | BallDroppings_Lite_1, Pro_Basketball_Scores_2 |
| 5 | game, score, jump | City_Jump_1, Jumper0 |
| 6 | score, pointf, high, ball, basket | BasketBall_v1, One_Touch_Battery_Saver*, Tranquilize_v1*, Balance_The_Beer_1 |
| 7 | view, game, flash, paddle | Adobe_Flash_Player_10*, Paddle_Bounce_1 |
| 8 | call, phone, phoneq, date, number | Jumper0*, PhoneQ_Lite_1 |
| 9 | camera, image, preview, effect, photo | Camera_Fun_Free_2, Camera_Illusion_1 |
| 10 | ad, message, air, game, logisoft | AirAttack_Lite_3, Air_Control_Lite_1 |
| 11 | weather, update, city, locate, temperature, wind, cloud | Animated_Weather_Free_2, CityWeather_1, Weather_notification_0 |

lists some automatically generated software categories that have more than two software applications. Some applications in the table are marked with "*" because they are misclassified. From the table we can find that the extracted topic words can effectively represent the themes of the categories. For example, applications in Category 1 and 3 are all related to multimedia, and applications in Category 4-7 and 10 are related to games. In the categorization of LACTA, these applications are classified into more specific categories.

## V. THREATS TO VALIDITY

Although LACTA shows its improvements in automatic software categorization for Android applications, there are some factors that may imperil the validity of the experimental results of LACTA. For threats of internal validity, one major concern is that many manual inspections are involved in the experiments to decide whether the extract topic words belong to some categories. Therefore, the subjectivity of judgments may be introduced. Since the development of Android applications is still emerging and the classifications for Android applications are very divergent in many software repositories, this problem cannot be avoided in the current situation. Another threat of internal validity is that the settings of thresholds are only studied for the collected applications. A more comprehensive study needs to be conducted to explore the generality of the effectiveness of LACTA. In addition, a threat of internal validity is that the performance of LACTA heavily relies on the exploration of Android domain knowledge. How to build a high quality knowledge base will be a major concern for the future study.

## VI. CONCLUSION

Noticing that the number of new Android applications tends to be rapidly increased in the near future, we find that automatic software categorization will be in great demand for management of Android application archives. Although there have been may approaches proposed to address the automatic software categorization problem, they cannot be directly applied to Android applications because most disseminated Android applications are distributed in the native Android application package (APK) from. To perform software categorization on Android applications needs the decompilation process which in turn complicates the automatic software categorization problem.

In this paper, we propose an enhancement based on the LACT approach [8] extensively employs Android domain knowledge with Latent Dirichlet Allocation (LDA) to improve the automatic software categorization performance for Android applications. In LACTA, software categories are first determined and then Android applications are classified into these categories accordingly.

Using 42 Android applications collected from the Internet, we conducted empirical experiments to evaluate the effectiveness of LACTA. Compared with the original LACT, LACTA shows its prominent improvements in the experiments. The promising performance shows the potential feasibility of LACTA.

There are still several issues which need to be discussed in the future. First, we plan to collect more Android applications for a comprehensive study of various application types. Second, the settings of thresholds used in LACTA play an important role in classification performance. A mechanism that can automatically determine these settings will facilitate automatic software categorization in practice. This challenging work will be also included in our plan. Finally, we will investigate the mechanisms to extract more meaningful features from Android applications. With more informative features, we are convinced that automatic software categorization for Android applications will be more feasible in daily use.

## REFERENCES

[1] Y. S. Maarek, D. M. Berry, and G. E. Kaiser, "An information retrieval approach for automatically constructing software libraries," *IEEE Trans. Softw. Eng.*, vol. 17, no. 8, pp. 800-813, Aug. 1991.

[2] D. Merkl, "Content-based software categorization by self-organization," in *Proc. of the IEEE International Conf. on Neural Networks*, 1995, pp. 1086–1091.

[3] A. Chan and T. Spracken, "Discovering common features in software code using self-organizing maps," in *Proc. of the International Symposium on Computational Intelligence (ISCI 2000)*, Kosice Slovakia, Aug. 2000.

[4] S. Ugurel, R. Krovetz, and C. L. Giles, "What's the code? Automatic classification of source code archives," in *Proc. of the 8th ACM SIGKDD International Conf. on Knowledge Discovery and Data Mining (KDD '02)*, 2002, pp. 632-638.

[5] R. Krovetz, S. Ugurel, and C. L. Giles, "Classification of source code archives," in *Proc. of the 26th Annual International ACM SIGIR Conf. on Research and Development in Information Retrieval (SIGIR '03)*, 2003, pp. 425-426.

[6] S. Kawaguchi, P. K. Garg, M. Matsushita, and K. Inoue, "MUDABlue: An automatic categorization system for open source repositories," in *Proc. of the 11th Asia-Pacific Software Engineering Conf.* (APSEC 2004), 2004, pp. 184–193.

[7] S. Kawaguchi, P. K. Garg, M. Matsushita, and K. Inoue. "MUDABlue: an automatic categorization system for open source repositories," *J. Systems and Software.*, vol. 79, no. 7, July 2006, pp. 939-953.

[8] K. Tian, M. Revelle, and D. Poshyvanyk, "Using latent Dirichlet allocation for automatic categorization of software," in *Proc. of the 6th IEEE International Working Conf. on Mining Software Repositories* (MSR 2009), 2009, pp.163-166.

[9] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent Dirichlet allocation," *Journal of Machine Learning Research*, vol. 3, 2003, pp. 993-1022.

[10] T. Landauer, P. W. Foltz, and D. Laham, "An introduction to latent semantic analysis," *Discourse Processes*, no. 25, 1998, pp. 259–284.

[11] A. Kuhn, S. Ducasse, and T. Gírba, "Semantic clustering: identifying topics in source code," *Information and Software Technology*, vol. 49, no. 3, 2007, pp.230-243.

[12] The dex2jar project. Available: http://code.google.com/p/dex2jar/

[13] The JD-GUI tool. Available: http://java.decompiler.free.fr/?q=jdgui

[14] The android-apktool project. Available: http://code.google.com/p/android-apktool/

[15] G. Maskeri, S. Sarkar, and K. Heafield, "Mining business topics in source code using latent Dirichlet allocation," in *Proc. of the First India Software Engineering Conf.*, Hyderabad, India, 2008, pp. 113-120.

[16] X.-H. Phan and C.-T. Nguyen, "GibbsLDA++: A C/C++ implementation of latent Dirichlet allocation". Available: http://gibbslda.sourceforge.net/