# Empirical Study of Co-Renamed Identifiers

Yuki Osumi, Naotaka Umekawa, Hitomi Komata, and Shinpei Hayashi

*Tokyo Institute of Technology*, Tokyo 152–8550, Japan

osumi@se.c.titech.ac.jp, hayashi@c.titech.ac.jp

*Abstract*—*Background:* **The renaming of program identifiers is the most common refactoring operation. Because some identifiers are related to each other, developers may need to rename related identifiers together.** *Aims:* **To understand how developers rename multiple identifiers simultaneously, it is necessary to consider the relationships between identifiers in the program and the brief matching for non-identical but semantically similar identifiers.** *Method:* **We investigate the relationships between co-renamed identifiers and identify the types of their relationships that contribute to improving the recommendation using more than 1M of renaming instances collected from the histories of open-source software projects. We also evaluate and compare the impact of co-renaming and the relationships between identifiers when inflections occur in the words in identifiers are taken into account.** *Results:* **We revealed several relationships of identifiers that are frequently found in the co-renamed identifiers, such as the identifiers of** *methods in the same class* **or** *an identifier defining a variable and another used for initializing the variable*, **depending on the type of the renamed identifiers. Additionally, the consideration of inflections did not affect the tendency of the relationships.** *Conclusion:* **These results suggest an approach that prioritizes the identifiers to be recommended depending on their types and the type of the renamed identifier.**

*Index Terms*—**identifier, rename refactoring, refactoring recommendation**

## I. INTRODUCTION

Program identifiers play a significant role in program comprehension, and identifiers are frequently renamed during software refactoring. It is estimated that identifiers account for approximately 70% of the total program content [1]. Corazza *et al.* also reported that developers can facilitate knowledge transfer among developers by naming identifiers appropriately that reflect their intentions and the domain knowledge [2]. Developers can infer the intent and the behavior of code fragments if the identifiers in them have appropriate names [3]. Therefore, appropriately naming identifiers and/or renaming them can make it easier for developers to guess the roles of identifiers and improve the understanding of the program. Renaming identifiers is considered as the most common form of refactoring [4].

When an identifier is renamed, it may be necessary to rename other identifiers related to the renamed identifier. This is because there may be a misunderstanding if multiple expressions refer to the same concept to be used for naming identifiers [1]. However, it is time-consuming to identify all the identifiers that need to be renamed. The replacement feature of source code editors and rename refactoring tools provided by current integrated development environments (IDEs) cannot determine all the identifiers that need to be renamed together. Therefore, developers must determine the identifiers to be renamed manually, which is difficult and may lead to missing changes.

In this study, we investigated co-renamed program identifiers to discuss their critical characteristics for improving the recommendation of identifiers to be co-renamed. We surveyed co-renamed identifiers extracted from past histories of open-source software (OSS) repositories, clarified the relationships among them that were prone to be correlated, and discussed the relationships that should be prioritized in the recommendation. In addition, we evaluated the variation of word stemming in identifiers and discussed the influence of inflections, i.e., the influence of singular and plural and those of participles, on the recommendation. There are several attempts to utilize the relationship among identifiers for renaming recommendations [5], [6]. The results of our empirical study can strengthen their approaches by prioritizing the relationships to be used for the renaming recommendation.

The main contributions of this work are summarized as follows:

- a large-scale empirical study using OSS repositories on the relationship between co-renamed identifiers and
- results of evaluating the effect of inflections in identifiers on co-renaming.

As a result of evaluating co-renamed identifiers, we identified the relationships among identifiers that frequently occur in renamed identifiers. Additionally, the frequency of relationships differed depending on the type of renamed identifiers. Furthermore, we found that the consideration of inflections did not affect the tendency of the relationships.

The rest of this paper is organized as follows. In the next section, we explain our aim to support identifier co-renaming via an example of co-renamed identifiers. Section III presents our evaluation methods. Section IV explains our evaluation results and discusses their significance for renaming recommendation. In Section V, we summarize existing work on identifiers and their renaming. Finally, Section VI concludes this paper and summarizes future work.

## II. MOTIVATION

When an identifier is renamed, another identifier may need to be renamed accordingly. This is because the existence of multiple representations of the same concept may cause developers to misunderstand the program [1]. An example of renamed identifiers is shown in Fig. 1[1]. In this commit, the class `MetricType` is renamed to `MetricAttribute`.

---

[1]https://github.com/dropwizard/metrics/commit/3ccd7a1

```
⋮
import com.codahale.metrics.MetricType;
⋮
public class GangliaReporter extends ScheduledReporter {
    ⋮
    private void announceIfEnabled(MetricType metricType, ...)
            throws GangliaException {
        if (getDisabledMetricTypes().contains(metricType)) {
            ⋮
    }
    ⋮
    private void announce(..., GMetricType type, ...)
            throws GangliaException {
        ⋮
```

(a) Before renaming.

```
⋮
import com.codahale.metrics.MetricAttribute;
⋮
public class GangliaReporter extends ScheduledReporter {
    ⋮
    private void announceIfEnabled(MetricAttribute metricAttribute, ...)
            throws GangliaException {
        if (getDisabledMetricAttributes().contains(metricAttribute)) {
            ⋮
    }
    ⋮
    private void announce(..., GMetricType type, ...)
            throws GangliaException {
        ⋮
```

(b) After renaming.

Fig. 1. Renamed identifiers in Metric project.

From the renaming, we can infer that the developer intended to replace the word "*type*" in the identifiers with "*attribute*". The gray-highlighted identifiers in the figure are those renamed in this commit. The blue-highlighted parts in the figure represent the changes based on the developer's intention, i.e., "*type*" before the renaming to "*attribute*" after the renaming. In contrast, the green parts in the figure indicate that they were not renamed; the word "*type*" remained after this co-renaming operations.

However, it is time-consuming for developers to determine all the identifiers that should be renamed. For example, although the source code editors' *Replace All* feature can replace all the matched strings simultaneously, it may lead to errors as it applies to elements that need not be replaced. Additionally, rename refactoring tools provided by IDEs such as Eclipse[2] and IntelliJ IDEA[3] can limit the identifiers to be renamed at once to those of the same program element; developers cannot change multiple identifiers to be renamed at once. Therefore, developers need to identify all the identifiers to be renamed simultaneously, which is a challenging task wasting much time.

According to the renamed one, recommending possible candidates of related identifiers to be renamed together can assist developers in improving the consistency of their source code. When a developer renames the class `MetricType` to `MetricAttribute` as shown in Fig. 1, by applying the substitution from the word "*type*" to "*attribute*" to the variable `metricType`, which is typed as `MetricType`, we

[2]https://www.eclipse.org/
[3]https://www.jetbrains.com/idea/

can recommend a renaming to `metricAttribute`.

In recommending identifiers to be renamed, the following should be taken into account:

- relationships between identifiers that are likely to be renamed together and
- inflections in identifiers.

Because not all identifiers should be renamed together to follow the developer's renaming intention, it is useful to consider the relationship between identifiers that are likely to be co-renamed. For example, on one hand, the substitution in Fig. 1 from the word "*type*" to "*attribute*" in the class `MetricType` is also applicable to the word "*type*" in the class `GMetricType`, but this did not happen; `GMetricType` was not renamed. On the other hand, this substitution was applied to the variable `metricType`, whose type is `MetricType`, and the variable identifier was renamed to `metricAttribute`. Here, a relationship of "variable and its type" exists between these two identifiers: the class `MetricType` and the variable `metricType`. It is useful to clarify the kinds of relationships between the identifiers that are likely to be renamed simultaneously.

Additionally, if inflections are not taken into account, the developer's renaming intention cannot be applied, and recommendations may fail. For example, in the renaming in Fig. 1, the method `getDisabledMetricTypes` has also been renamed, and a similar substitution from the word "*type*" to "*attribute*" has occurred; however, the plural form "*types*" has been replaced with "*attributes*", unlike the case of `MetricType`. These renamings should be considered as those based on the same intention.

To clarify the characteristics of the relationships between identifiers to be co-renamed together and the effects of inflections in identifiers, we conducted an empirical study of co-renamed identifiers. Through answers to the following research questions (RQs), we discuss the aspects that should be considered in recommending identifiers to be co-renamed.

- $RQ_1$: How often do co-renamings occur?
- $RQ_2$: To what extent do co-renamed identifiers correlate the relationships of identifiers?
- $RQ_3$: What is the difference when the inflections in identifiers are taken into account?

In $RQ_1$, we confirm that a large number of co-renamed identifiers occur. In $RQ_2$, we characterize the relationships between identifiers that have undergone co-renamings. In $RQ_3$, we clarify the effects of inflections in identifiers on the number of co-renamings and the relationships between identifiers that have undergone co-renamings.

## III. METHODOLOGY

### A. Overview

Figure 2 shows the overview of our study. First, we extract all the renamings from the commit history of the version control repositories using an existing refactoring detection tool. Next, we lemmatize each word in the identifier names before and after each extracted renaming to ignore the inflection.
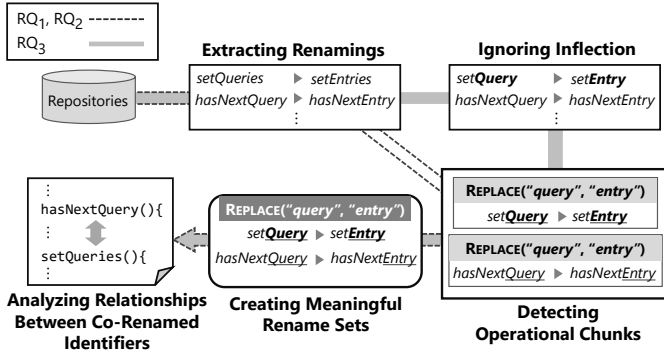
Fig. 2. Overview of our study.

Subsequently, we compare the old identifier name with the new one to obtain the operational chunks of the renaming. We then create meaningful rename sets that share their operational chunks. Finally, we examine the relationship between their target identifier names in the source codes for each pair of two renamings in a meaningful rename set. Note that we skip the third step when answering $RQ_1$ and $RQ_2$ since we consider the inflection only when answering $RQ_3$.

### B. Extracting Renamings

We detected renamings using RefactoringMiner ver. 2.0.2 [7], [8] in 187 of the repositories used by Silva *et al.* and Tsantalis *et al.* in their refactoring studies [7], [9]. For each repository, we extracted class renamings (Rename Class), method renamings (Rename Method), class attribute renamings (Rename Attribute), method parameter renamings (Rename Parameter), and variable renamings (Rename Variable) from detected refactorings and created a set $\mathbb{R}$ of all the renamings. Let $r.commit$ denote the commit in which a renaming $r \in \mathbb{R}$ was performed. We excluded repositories that did not contain any renamings, resulting in 176 repositories with a history from June 2001 to January 2021. The total number of commits in these repositories was 1,883,276, whereas the total number of extracted renamings was 1,084,121.

### C. Ignoring Inflection

We split identifier names into words using Ronin in the Spiral [10] package for each identifier name before and after a renaming. For each word, we converted it to lower case and lemmatized it to ignore inflection using WordNetLemmatizer [11].

### D. Detecting Operational Chunks

Using a difference extraction algorithm based on the longest continuous matching subsequence detection, we extracted each consecutive addition, deletion, and replacement of words as one *operational chunk* from the identifier names before and after a renaming $r$. We denote the all set of operational chunks of $r$ by $r.chunks$. For each renaming $r$, we split identifier names before and after the renaming $r$ into words using Ronin and created word sequences. We detected $r.chunks$ from the differences between the word sequence before and after

renaming $r$. An operational chunk consists of the sequence of words changed by renaming and changing types of words. By introducing operational chunks, we can regard renamings for different identifier names as a co-renaming that shares the same operational chunk.

There are four types of operational chunks as listed below.

- INSERT(*added words*): Addition of words. For example, a renaming `dataProviderId` → `dataProviderInstanceId` involves an operational chunk INSERT("*instance*").
- DELETE(*deleted words*): Deletion of words. For example, a renaming `SkipConstantResult` → `SkipResult` involves an operational chunk DELETE("*Constant*").
- REPLACE(*deleted words*, *added words*): Replace of words. For example, a renaming `getRandom` → `createRandom` involves an operational chunk REPLACE("*get*", "*create*").
- OTHER(*words*): Other type of renaming. For example, a renaming `TIMES` → `times` downcases its word, which belongs to an operational chunk OTHER("*time*").

It should be noted that OTHER is only detected when none of INSERT, DELETE, or REPLACE is detected in a given renaming. When detecting OTHER, we compared the identifier names before and after a renaming, starting from the first word, and detected OTHER for each word that differs.

### E. Creating Meaningful Rename Sets

Based on the operational chunks detected in Section III-D, we created a meaningful rename set, a collection of renamings that share the same operational chunk in the same commit. We defined a meaningful rename set $U_{c,h}$ in commit $c$, where its belonging operational chunk is $h$, as follows:

$$U_{c,h} := \{r \in \mathbb{R} \mid r.commit = c \;\wedge\; r.chunks \supseteq h\}.$$

For the set $\mathbb{R}$ of all the renamings in each repository, we created a set $\mathbb{U}$ of meaningful rename sets, i.e., $\mathbb{U} = \{U_{c,h} \mid U_{c,h} \neq \emptyset\}$.

When we detected multiple operational chunks in a renaming, we created meaningful rename sets for each operational chunk and included the renaming in all their sets. For example, renaming $r$: `minimumVersion` → `versionSpec` is an element of both $U_{\text{DELETE}("minimum")}$ and $U_{\text{INSERT}("spec")}$.

### F. Analyzing Relationships between Co-Renamed Identifiers

For each pair of two renamings in a meaningful rename set $U_{c,h}$, we evaluated the relationship between the renamed identifiers in the source codes. Let $P(U_{c,h})$ denote a set of all pairs of two renamings in $U_{c,h}$, i.e.,

$$P(U_{c,h}) := \{(r_i, r_j) \mid r_i, r_j \in U_{c,h} \;\wedge\; r_i \neq r_j\}.$$

For a pair $p_i = (r_i, r_j) \in P(U_{c,h})$, we detected relationships $R_c(p_i)$ from the source codes of the parent commit of commit $c$, i.e., the source codes before $r_i$ and $r_j$ perform as follows:

$$R_c(p_i) = relation_c(r_i.old, \; r_j.old)$$

where $r.old$ is the identifier before $r$ is performed and $relation_c(I_i, I_j)$ is a set of relationships between identifiers $I_i$ and $I_j$ in the source codes of the parent commit of $c$. We defined 14 relationships between co-renamed identifiers and categorized them into three: Location, Type, and Call/Data Dependency. The details of each relationship **R** are as follows.

1. Location

- BELONGS$_C$: Relationship between a class $c$ and its inner class $c.c$.
- BELONGS$_M$: Relationship between a class $c$ and its method $c.m$.
- BELONGS$_F$: Relationship between a class $c$ and its attribute $c.a$.
- BELONGS$_A$: Relationship between a method $m$ and its parameter $a$.
- BELONGS$_L$: Relationship between a method $m$ and its local variable $v$.
- CO-OCCURS$_M$: Relationship between two methods $c.m_1$ and $c.m_2$ in the same class $c$.

2. Type

- EXTENDS: Relationship between a class $c$ and its direct subclass $c'$.
- IMPLEMENTS: Relationship between an interface $i$ and the class $c$ that implements $i$.
- TYPE$_M$: Relationship between a method $m$ and its return type $c$.
- TYPE$_V$: Relationship between an attribute, a variable, or a method parameter $v$ and its type $c$.

3. Call and Data Dependency

- INVOKES: Relationship between a method $m_1$ and another method $m_2$ invoked by $m_1$ ($m_1 \neq m_2$).
- ACCESSES: Relationship between a method $c.m$ of a class $c$ and the attribute $c.a$ of the same class $c$ referenced by $c.m$.
- ASSIGNS: Relationship between a left side attribute or variable $v_{left}$ of an assignment statement and a right side attribute of the assignment statement $v_{right}$ such as an attribute, a parameter, a variable or an invocation of a method.
- PASSES: Relationship between a (formal) parameter $p$ of a method $m$ and an argument (actual parameter) $a$ of the method $m$ such as an attribute, a variable, or invocation of a method.

See again Fig. 1(a) with assuming that the developer triggered to rename the class name `MetricType` to `MetricAttribute`. In this code fragment, the type of the parameter `metricType` is `MetricType`, which means that there is a relationship of TYPE$_V$ between the parameter `metricType` and the class `MetricType`. Also, the method `getDisabledMetricTypes` returns an object typed as `Set<MetricType>`, which is not expressed in Fig. 1(a) though, which leads to a relationship of TYPE$_M$ between the method `getDisabledMetricTypes` and the class `MetricType`. We identify these kinds of relationships by applying simple static analyses for parsed source files.

We converted the target source codes to XML files using srcML [12] and detected $relation_c(r_i.old, r_j.old)$ using XPath. For each relationship **R**, we defined XPaths that return nodes only if a pair $(r_i.old, r_j.old)$ satisfies **R**. For example, in BELONGS$_M$, we defined an XPath that returns nodes only if a class $c$ has a method $m$ as follows.

- //class[name[text()=$c$]]/block/function[name[text()=$m$]]

For a pair $(r_i.old, r_j.old)$, we created a set of relationships **R** that the pair satisfied and regarded the set as a $relation_c(r_i.old, r_j.old)$.

## IV. EMPIRICAL STUDY

### A. $RQ_1$: How often do co-renamings occur?

*1) Motivation:* The purpose of this RQ is to confirm that many renamings co-occur. We determined the number of co-renamed identifiers and the characteristics of these renamings.

*2) Study Design:* For each of the 176 repositories in Section III-B, we evaluated a percentage of co-renamings. For each $\mathbb{U}$ created in each repository, we computed the ratio of the total number of elements of meaningful rename sets that contained more than one element $\sum_{U \in \mathbb{U}, |U| \geq 2} |U|$ to the total number of all elements of meaningful rename sets $\sum_{U \in \mathbb{U}} |U|$. We regarded this ratio as the rate of co-renamings.

We determined the ratio of the number of co-renamed identifiers to the number of all the co-renamed identifiers for each co-renamed identifier. For the case where $n$ identifiers are co-renamed, we computed its ratio of the total number $\sum_{U \in \mathbb{U}, |U|=n} |U|$ to all the co-renamings $\sum_{U \in \mathbb{U}, |U| \geq 2} |U|$.

For each number of co-renamed identifiers, we also determined the ratio of the number of unique identifier names included in the co-renamings to the number of co-renamed identifiers. Because some co-renames contain renamings for the same identifier names, the number of renamed identifier names in the co-renamings may be small even when the number of co-renamed identifiers is large. We computed the rate of the total number of meaningful rename sets in which $m$ unique identifier names were involved to the number of meaningful rename sets whose size was $n$. Here, the number of unique identifier names $m$ was computed as the number of non-overlapping pre-renaming identifier names $|\{r.old \mid r \in U\}|$ in a meaningful rename set.

*3) Results:* On average, 57% of the renamings were co-renamings. The left side of Fig. 3 is a box plot showing the ratio of the total number of elements in the meaningful rename sets with more than one element to the total number of elements in all the meaningful rename sets, i.e., the rate of co-renamings in each repository. The vertical axis represents the rate of co-renamings. Although the rates are widely distributed (16% to 89%), the first quartile is 51%, which shows that 75% of all projects contain more than 50% of co-renamings. The rate of co-renamings is 0.57 for both the mean and median.

Figure 4 shows a cumulative distribution of the percentage of the total number of elements of the meaningful rename sets for each number of the elements of the meaningful rename sets, colored by the number of unique identifier names. The
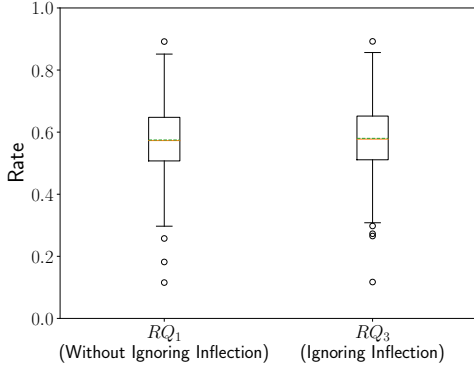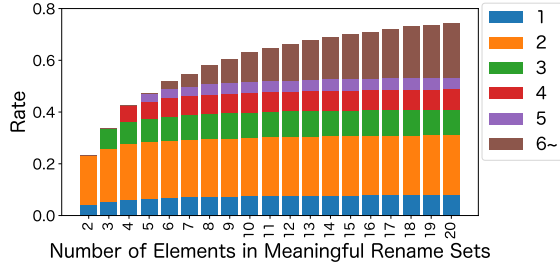
Fig. 3. Rate of co-renamings.



Fig. 4. Cumulative distribution of the rate of co-renamings, colored by the number of unique identifier names.



Fig. 5. Rate of relationships.

vertical axis represents the rate of meaningful rename sets, whereas the horizontal axis represents the maximum number of elements of the meaningful rename sets. For example, the blue area in the figure shows the rate of meaningful rename sets that have only one unique identifier name.

The number of elements of a co-renaming was small. From Fig. 4, we can observe that meaningful rename sets with less than six elements account for more than 50% of all the meaningful rename sets considered as co-renaming. Furthermore, 70% of all the meaningful rename sets had no more than 20 elements.

On the other hand, most of the co-renamings consisted of different identifier names. Figure 4 shows that the rate of co-renamings with less than five unique identifier names did not increase even if the number of elements in meaningful rename sets increased. In other words, there were few identifiers in a meaningful rename set whose names were the same as those of other identifiers in the meaningful rename set. This result indicates that when recommending co-renamings, it is essential to consider how identifier names were renamed, i.e., operational chunks, not only to recommend identifiers matching names with those of renamed identifiers.

> 57% of all the renamings were co-renamings. The co-renamings were composed of relatively few renamings. It is important to consider operational chunks when recommending co-renamings.
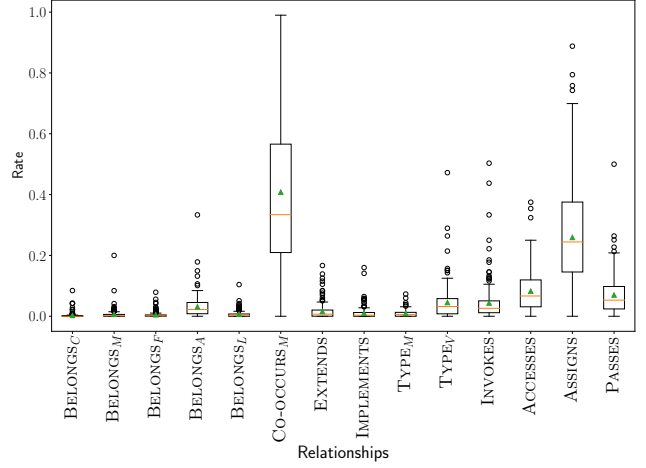
### B. $RQ_2$: To what extent do co-renamed identifiers correlate the relationships of identifiers?

*1) Motivation:* Because it is unnecessary to rename all the other identifiers to which an operational chunk of a renaming can be applied, we reveal the relationships between identifiers that are likely to be co-renamed. For co-renamings detected in each repository, we evaluated the relationship between the co-renamed identifiers.

*2) Study Design:* We determined the relationships between the co-renamed identifiers and evaluated the 176 repositories, same as $RQ_1$. From these repositories, 653,194 meaningful rename sets were created. Of these, we analyzed 138,250 meaningful rename sets that contained more than two elements and detected 732,390 relationships.

We also limited the detection of relationships to meaningful rename sets that contain at least one renaming for a specific type of identifiers: Class, Method, Attribute, Parameter, and Variable. By limiting the analysis, we could detect relationships between these types of identifiers, and those that are likely to be co-renamed. For example, if we limit the detection of relationships to a meaningful rename set that contains at least one renaming of a Class, we do not analyze the meaningful rename sets that do not contain renamings of Classes.

*3) Results:* The most detected relationships were $\textsc{Co-occurs}_M$ and $\textsc{Assigns}$. Figure 5 shows a box plot of the rate of each relationship. The vertical axis represents the rate of relationships, whereas the horizontal axis represents the kinds of relationships. The mean values for the rate of $\textsc{Co-occurs}_M$ and $\textsc{Assigns}$ were 40.8% and 25.9%, respectively, the highest of all the relationships, while the other relationships were much lower.

Figure 6 shows a box plot of the rate of each relationship, limiting the analysis to co-renamings containing a specific type. In each figure, the vertical axis represents the rate of relationships, whereas the horizontal axis represents the kinds of relationships. Each box plot shows the distribution

(a) Classes.　　　　(b) Methods.　　　　(c) Attributes.
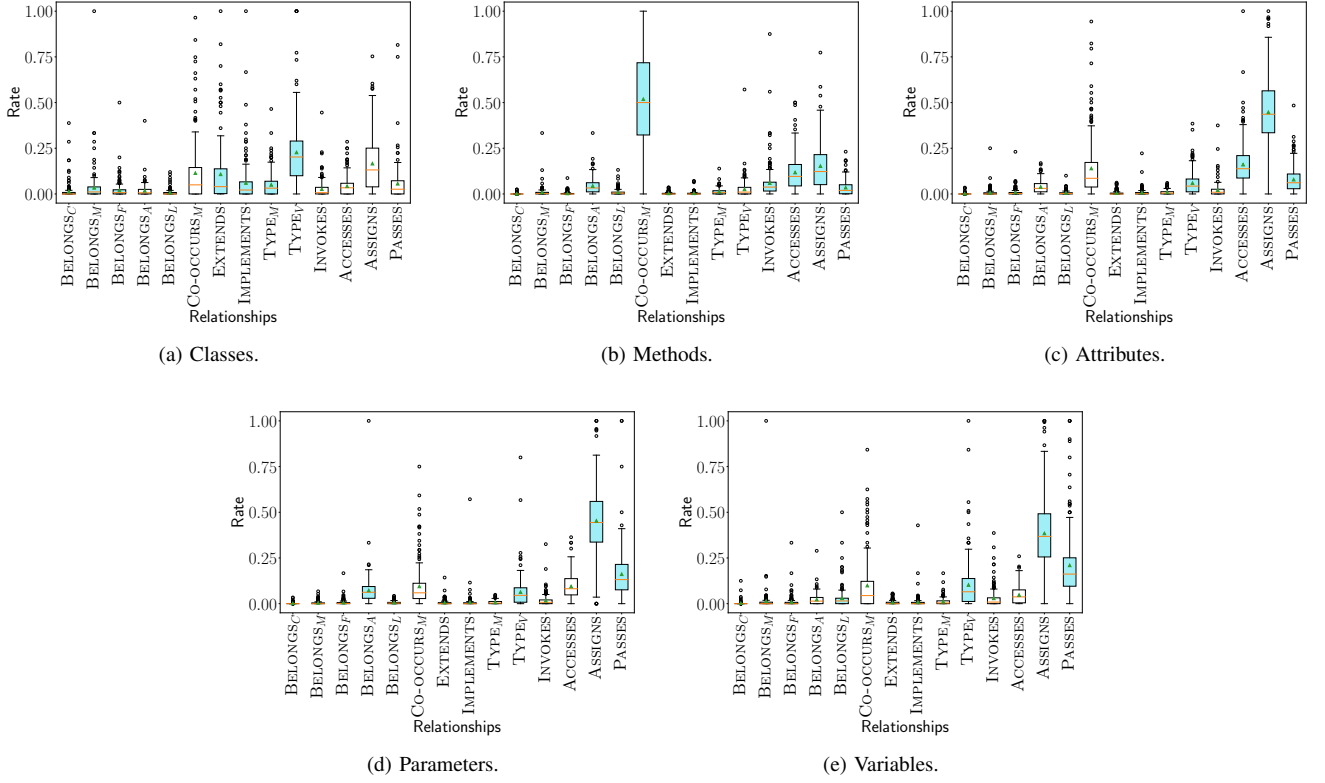


(d) Parameters.　　　　(e) Variables.

Fig. 6. Rate of relationships detected in co-renamings containing renamings of specific identifiers.

of particular relationships detected to the total relationships detected for each repository. Light blue-colored ones are for the relationships that involve limited types of identifiers. For example, Fig. 6a shows the rate of relationships when limiting the analysis to co-renamings that contain at least one renaming of Class, and the box plots of relationships that involve class renamings, i.e., BELONGS$_C$, BELONGS$_M$, BELONGS$_F$, EXTENDS, IMPLEMENTS, TYPE$_M$, and TYPE$_V$ are highlighted as light blue.

The different distributions for each type of identifier in Fig. 6 indicate that the relationships in which identifiers are likely to be co-renamed differ depending on the type of identifier renamed. CO-OCCURS$_M$ is the most detected relationship for Method, ASSIGNS for the other types. Note that TYPE$_V$ is frequently detected relationship in Class; ACCESSES in Attribute, PASSES in Parameters and Variables. Each of these relationships is associated with the corresponding identifier type. This result means that identifiers in a relationship to a renamed identifier are likely to be co-renamed with the identifier.

> The distribution of the relationships differed significantly among the types of identifiers. The most detected relationship was CO-OCCURS$_M$ for Method, ASSIGNS for the other types. In addition, TYPE$_V$ frequently appeared in Class; ACCESSES in Attribute, PASSES in Parameter and Variable. These results suggest that identifiers that are likely to be co-renamed depend on the type of identifiers.

### C. RQ$_3$: What is the difference when the inflections in identifiers are taken into account?

*1) Motivation:* In $RQ_1$ and $RQ_2$, we did not ignore inflection. Thus, operational chunks that should be the same in multiple renaming may become a different operational chunk in each renaming. Therefore, we evaluated the rate of co-renamings and the tendency of the relationships between the renamed identifiers in both cases of ignoring inflection and not and clarifying the effect of the inflection.

*2) Study Design:* In order to take inflections into account, when detecting operational chunks (in Section III-D), we ignored inflection. After splitting identifier names into word sequences in Section III-D, for each word, we lemmatized it to ignore inflection using WordNetLemmatizer [11].

WordNetLemmatizer can convert word forms within a part of speech, such as the singular and plural forms of nouns and the present and past tenses of verbs, but cannot convert word forms across parts of speech. For example, the plural form of a noun such as *queries* can be converted to its singular form *query*, but the noun *creator* cannot be converted to the verb *create*.

Then, we evaluated the effect of inflection on the operational chunks of renamings. We evaluated the 176 repositories used in $RQ_1$. If we do not ignore inflection, the operational chunks of renaming change only the inflection REPLACE instead of OTHER. For example, in a renaming node → nodes, when ignoring inflection an operational chunk OTHER("*node*") is detected, while not an operational chunk REPLACE("*node*",
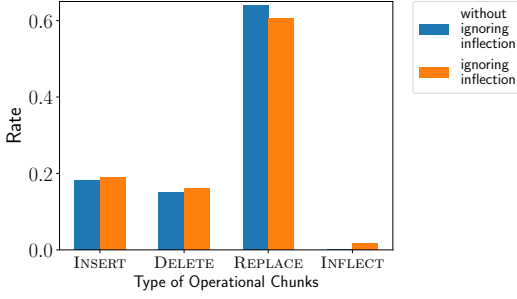
Fig. 7. Rate of detected operational chunks.



Fig. 8. Rate of relationships in co-renamings affected by inflection.

"*nodes*") is detected. To evaluate the difference in detecting such operational chunks, we defined a new operational chunk INFLECT, which represents inflection.

- INFLECT(*words*): Inflection of words. For example, a renaming of `instance` → `instances` involves an operational chunk of INFLECT("*instance*").

When we detected OTHER from renaming $r$, we compared the words in the identifiers before and after the renaming from the beginning in a case-insensitive manner and detected INFLECT each word that differs. We compared the rate of each operational chunk, including INFLECT, ignoring and not ignoring inflection.

We also evaluated the effect of inflection on the co-renamings in the same 176 repositories as $RQ_1$ and $RQ_2$. We evaluated the difference in the ratio of the total number of co-renamings to the total number of renamings and that in the number of meaningful rename sets $|\mathbb{U}|$ between the two cases of ignoring and not ignoring inflection for each repository.

Additionally, We evaluated the effect of the relationships between co-renamed identifiers in the same 176 repositories as in $RQ_2$. For each repository, we created a set of meaningful rename sets $\mathbb{U} = \{U_1, U_2, \ldots\}$ ignoring inflection and a set of meaningful rename sets $\mathbb{U}^{skip} = \{U_1^{skip}, U_2^{skip}, \ldots\}$ without ignoring inflection. We analyzed relationships between identifiers in the newly created meaningful rename sets, ignoring inflection, i.e., the meaningful rename sets included in a relative complement of $\mathbb{U}^{skip}$ in $\mathbb{U}$ ($\mathbb{U} \setminus \mathbb{U}^{skip}$). The total number of newly created meaningful rename sets in all repositories was 119,165.

*3) Results:* Figure 7 shows the results of the evaluation of differences in operational chunks with and without ignoring inflection. Figure 7 shows the rate of operational chunks detected for each case with and without ignoring inflection in all repositories. The vertical axis shows the rate of the total number of operational chunks, and the horizontal axis shows the types of operational chunks. We excluded OTHER from the graph to focus on whether operational chunks changed due to inflection. The rate of OTHER was 0.03 in both cases.

The rate of the number of operational chunks did not change significantly owing to the effect of inflection. Figure 7 shows that when we ignored inflection, the rate of REPLACE decreased, whereas the rate of INFLECT increased. One of the reasons for this result is that the operational chunk REPLACE,
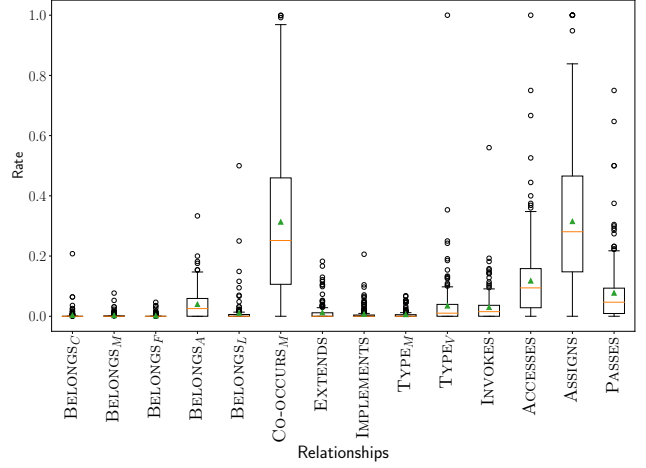
detected without ignoring inflection, changes to the operational chunk INFLECT due to ignoring inflection. However, the decrease in the rate of REPLACE was slight and did not significantly change the overall distribution.

The effect of inflection on the ratio of co-renamings to the total renaming was negligible. The right side of Fig. 3 shows the rate of co-renamings detected by ignoring inflection. The distribution, whose mean value of the rate of co-renamings detected with ignoring inflection is 0.58, is almost the same as that without ignoring inflection, as shown in the left side of Fig. 3 used in answering $RQ_1$. The total number of meaningful rename sets when ignoring inflection was 651,389, whereas that when not ignoring inflection was 653,194, with almost no difference between the two cases. However, there were cases where the number of detected operational chunks changed owing to ignoring inflection, and the number of meaningful rename sets also changed. Because a meaningful rename set is created for each operational chunk, renaming that increases the number of operational chunks by ignoring inflection may create a new meaningful rename set and affect the results.

Figures 8 and 9 show the results of analyzing relationships for the co-renamings affected by inflection. Figure 8 is a box plot of the results of the rate of each relationship for each repository, as in Fig. 5; Fig. 9 is a box plot of the results of the rate of each relationship obtained by limiting the detection of relations to co-renamings containing renamings of the specific type of identifiers for each repository, as in Fig. 6.

The rate of relationships did not change significantly in general, except when analyzing co-renamings containing renamings of Class. Comparing Fig. 8 with Fig. 5, the rate of each relationship somewhat changed but not considerably different. Comparing Fig. 9 with Fig. 6, the rate of each relationship is also not considerably different for Method, Parameter, and Variable; however, for Class, the rate of $\text{TYPE}_V$ increases, while many other rates decrease. This result may be due to ignoring the word inflection of the identifier name. For example, if the instances *query* and *queries* of the class *Query*
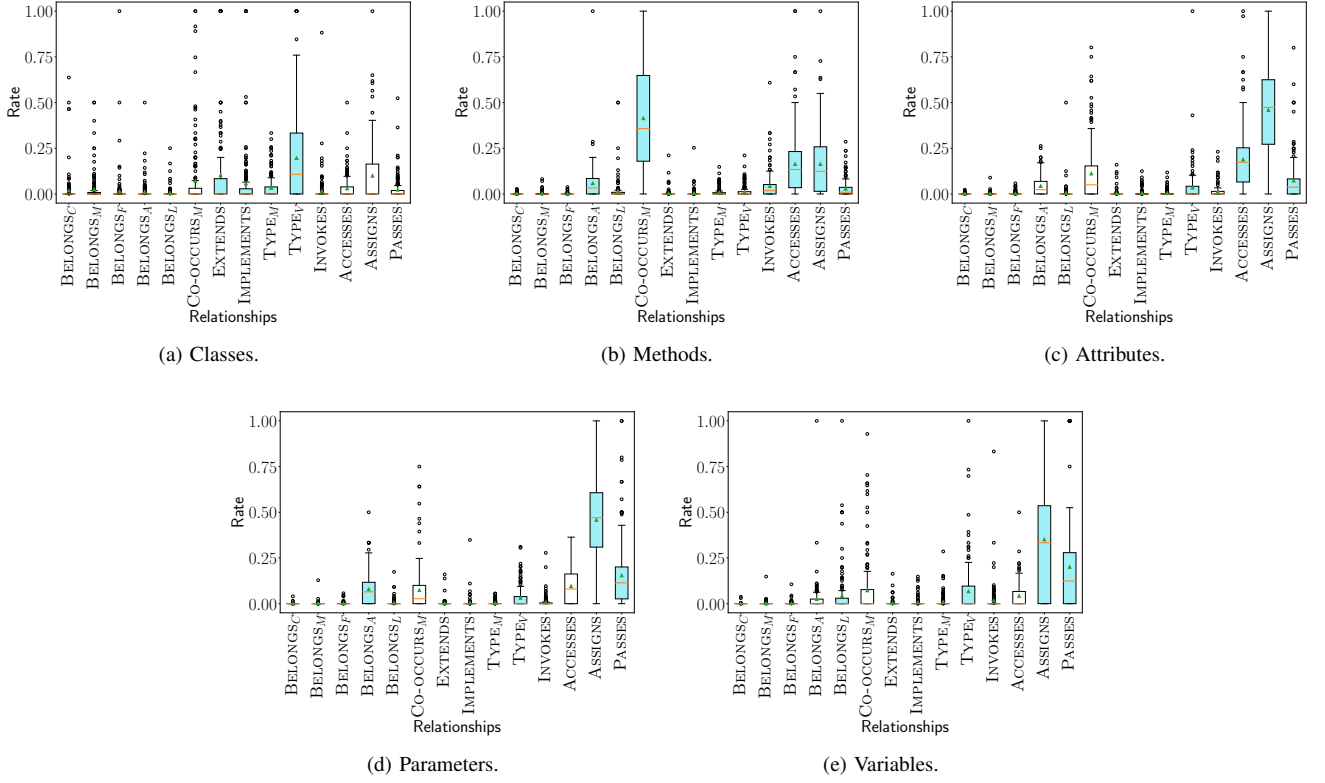
Fig. 9. Rate of detected relationships in co-renamings affected by inflection and containing renamings of specific identifiers.

are co-renamed to *entry*, *entries*, and *Entry*, respectively, the number of $\text{TYPE}_V$ increases because of ignoring the inflection of the identifier name. We assume that the rate of $\text{TYPE}_V$ increased owing to such co-renamings.

> There was a slight effect of inflection on the rate of operational chunks. It also had no significant effect on the rate of co-renamings. There was no significant effect in terms of relationships in general, except for an increase in the rate of $\text{TYPE}_V$ in the co-renamings containing renamings of Class.

### D. Discussion

In recommending co-renamings, we should change the identifiers recommended preferentially based on the type of renamed identifier and relationships between them. In $RQ_2$, we found that the tendency of relationships between renamed identifiers was different for each type of renamed identifier. We inferred that we can improve the accuracy of recommending co-renamings by recommending identifiers that have a relationship, which are likely to be co-renamed with a renamed identifier, based on the type of renamed identifier. For example, suppose that a class `Sample` has methods `addItem` and `removeItem`, and `addItem` is renamed to `addElement`. This renaming is for a method, and it involves REPLACE("*item*", "*element*"). In answering $RQ_2$, we found that the most detected relationship is CO-OCCURS$_M$ in method renamings. Therefore, we can recommend renaming identifiers in a CO-OCCURS$_M$ relationship with `addItem`,

i.e., `removeItem` to `removeElement`. We found a case similar to this example in neo4j[4]. In this project, there were renamings including DELETE("*count*") such as `countPins` → `pins` and `countBytesRead` → `bytesRead`. These renamings were only for methods; no renamings occurred for other type identifiers including "*count*" like an attribute `pageCount`. Changes in interface methods probably caused this, so it was not necessary to rename all the identifiers, including "*count*".

However, the effect of inflection is negligible, and we should pay little attention to the inflection. In $RQ_3$, we evaluated the changes in the rate of operational chunks and co-renamings due to ignoring inflection and the trend in relationships between co-renamed identifiers in meaningful rename sets affected by ignoring inflection. However, the results of $RQ_3$ did not significantly differ from the results of $RQ_1$ and $RQ_2$. This result suggests that we do not need to pay attention to inflection in recommending co-renamings. Although the detection rate of $\text{TYPE}_V$ increased for the co-renamings containing renamings of Class when ignoring inflection, this relation also had a high detection rate when not ignoring inflection. Therefore, inflection does not seem to affect the recommendation of co-renamings. However, in some cases, the number and type of detected operative chunks changed depending on whether we ignored inflection or not. Thus, it is necessary to evaluate the case such that the number of

[4]https://github.com/neo4j/neo4j/commit/f24af6d

operational chunks changes because of ignoring inflection in detail.

Some relationships may have a direction of propagation of renamings, and it is necessary to evaluate the directions of the relationships. We can consider the direction of a relationship when one renaming of two renamings in a relationship causes the other renaming. In $RQ_3$, we detected many PASSES for co-renamings containing renamings of Parameters or renamings of Variables. Liu *et al.* [13] evaluated identifiers with PASSES relationships and found that in many cases, it is better to rename arguments (i.e., Variables) than Parameters. Based on Liu *et al.* and our results, we deduced that most of the co-renamings in the PASSES relationship are due to renamings of Parameters, which causes renamings of Variables. By evaluating the directions of relationships, we may discover more valuable things for a recommendation of co-renamings.

*E. Threats to Validity*

*1) Internal Validity:* The tools used in this study may have influenced the results. We used RefactoringMiner [7] to extract renamings; but it is not the latest version, but the October 2018 version. However, the latest version of RefactoringMiner might yield different results from our evaluation. In addition, we used WordNetLemmatizer [11] to ignore inflection, but we did not care about the accuracy of the tool.

The relationships between identifiers in our evaluation do not cover all the relationships between the identifiers that tend to be co-renamed. In fact, in several meaningful rename sets, we could not detect any relationships. We can evaluate the relationships in more detail by examining these meaningful rename sets and defining new relationships.

*2) External Validity:* The results of our evaluation may not be the same for all repositories in general. We evaluated only open-source Java repositories. In repositories we did not evaluate, the trend of co-renamings and the degree of the influence of inflection might be different from our results. We evaluated 176 repositories of various sizes to address this issue: from 2,482 to 1,907,462 LOC and from 131 to 255,628 commits.

## V. RELATED WORK

Identifiers in a program have an important role in program understanding [14]–[16]. Deissenboeck and Pizka stated that identifiers account for approximately 70% of program texts [17]. Corazza *et al.* reported that developers could smoothly communicate with each other by giving appropriate identifier names that reflect the developer intention and the domain knowledge [2]. If identifiers are named appropriately, developers can infer their intention and behavior [3].

Peruma examined identifier renamings along with commit messages, their data types, and refactorings before and after them and reported that developers tended to rename identifiers to narrow their meanings, and 17.39% of all the renamings changed their corresponding data types [18].

Arnaoudova *et al.* proposed REPENT [19], an approach to automatically detect and classify identifier renamings in source code. Based on a natural language processing technique, REPENT classifies renamings into different categories such as meaning-preserved, narrowed, or broadened. Peruma *et al.* presented an empirical study of how identifiers were renamed, with an attention of whether the meaning of identifiers were to be narrowed or broadened [20]. Their classifications are more semantic-oriented, whereas our types of operational chunks are more lexical, based on the sequence of words in identifiers. In addition, REPENT and the study by Peruma *et al.* is focused on the classification of renaming instances, whereas our analysis focused on the co-renamings, i.e., the relationship between renamed identifiers.

Several studies tried to perform a parsing for investigating the structure of identifiers with a specialized grammar and capturing a better meaning of the word fragments in identifiers [21], [22]. In addition, abbreviations are common in identifiers, and the expansion of such abbreviations can reveal the meaning of identifiers and relationships among them [23]–[25]. These approaches can be more reliable than a simple inflection detection in our approach, and an extension of the inflection analysis in this paper may have room for improvements by embedding such sophisticated analyses approaches into our study framework.

Techniques to correct identifier names have been proposed for avoiding inconsistent identifier names that hinder program understanding. Some of them normalize identifiers so that the naming conventions are consistent for the entire program. Caprile and Tonella proposed a technique to generate a new identifier name by normalization using pre-defined rules and dictionaries [15]. Surafel *et al.* proposed a technique to suggest identifier names using an ontology generated from source code [26]. Kashiwabara *et al.* proposed a method to recommend an appropriate verb to be used as a method name using association rule mining [27]. Lawrie *et al.* proposed a technique to calculate the similarity of vocabulary in identifiers using information retrieval technology and normalizing them [16].

Meanwhile, the renaming by developers frequently occurs during software refactoring. Because an inconsistent renaming leads to inconsistent identifier names, preventing such renaming is also important. Fowler's refactoring catalog [28] contains rename refactorings such as Rename Method. Renaming is the most used refactoring operation [4], [29], and developers are able to rename a specific identifier using existing refactoring tools consistently. However, these tools do not have a co-renaming feature of related identifiers. Liu *et al.* [5] proposed RenameExpander, which recommends the renamings of identifiers related to the identifier renamed by developers. However, they did not investigate what relationships frequently occurred. Thies and Roth proposed a rename refactoring recommendation approach based on the assignments for variables [6], which is similar to the relationship of ASSIGNS in our approach. Our results can provide not only empirical evidence of the effectiveness of such approaches but also supportive opportunities to improve their recommendation results with a deep look at the type of relationships.

Multiple renamings tend to occur at once in software refactoring. If there is an identifier related to another identifier to be renamed, it may be necessary to rename it too. Saika *et al.* analyzed a refactoring operation history performed by developers and reported that a series of refactorings such as multiple renamings is often performed [30]. There are several attempt of an IDE support for applying a sequence of continued refactoring operations [31]. The results of our empirical study can support the design of such refactoring tools. Programmer-friendly refactoring tools can be improved by being aware of the type of the identifier to be triggered to rename and providing tailored support according to what entities are to be renamed.

## VI. Conclusion

In this study, we evaluated the effects of co-renamed identifiers on the relationships between identifiers and on the word form changes of words in the identifiers. The results showed that half of the identifier renamings occurred together with other renamings. Additionally, the relationships between identifiers that are likely to be co-renamed differ depending on the type of the renamed identifiers. Finally, a slight effect was observed for word form changes. These results suggest that, in recommending a renaming, it is beneficial to prioritize candidate identifiers to be co-renamed for each type of the renamed identifier.

Future work will include further definition and evaluation of relationships, and the evaluation of the direction of relationships. In some co-renamings of identifiers, the relationships considered in this study were not detected at all. Defining new relationships for these co-renamings may lead to discovering additional characteristics for recommending new identifiers to be renamed. Additionally, there was a possibility that some of the relationships had a direction. It is possible that more detailed characteristics of the relationships can be revealed by determining the identifier renamings that cause another, e.g., by considering the time series of the changes.

The dataset of co-renamed identifiers used in this study is publicly available [32].

## References

[1] F. Deissenboeck and M. Pizka, "Concise and consistent naming," *Software Quality Journal*, vol. 14, no. 3, pp. 261–282, 2006.

[2] A. Corazza, S. D. Martino, and V. Maggio, "LINSEN: An efficient approach to split identifiers and expand abbreviations," in *Proc. ICSM*, 2012, pp. 233–242.

[3] W. C. Wake, *Refactoring Workbook*. Addison-Wesley, 2003.

[4] G. C. Murphy, M. Kersten, and L. Findlater, "How are Java software developers using the Eclipse IDE?" *IEEE Software*, vol. 23, no. 4, pp. 76–83, 2006.

[5] H. Liu, Q. Liu, Y. Liu, and Z. Wang, "Identifying renaming opportunities by expanding conducted rename refactorings," *IEEE Transactions on Software Engineering*, vol. 41, no. 9, pp. 887–900, 2015.

[6] A. Thies and C. Roth, "Recommending rename refactorings," in *Proc. RSSE*, 2010, pp. 1–5.

[7] N. Tsantalis, M. Mansouri, L. Eshkevari, D. Mazinanian, and D. Dig, "Accurate and efficient refactoring detection in commit history," in *Proc. ICSE*, 2018, pp. 483–494.

[8] N. Tsantalis, A. Ketkar, and D. Dig, "RefactoringMiner 2.0," *IEEE Transactions on Software Engineering*, vol. 48, no. 3, pp. 930–950, 2022.

[9] D. Silva, N. Tsantalis, and M. T. Valente, "Why we refactor? Confessions of GitHub contributors," in *Proc. FSE*, 2016, pp. 858–870.

[10] M. Hucka, "Spiral: Splitters for identifiers in source code files," *Journal of Open Source Software*, vol. 3, no. 24, 653, pp. 1–3, 2018.

[11] S. Bird, E. Klein, and E. Loper, *Natural language processing with Python: Analyzing text with the natural language toolkit.* O'Reilly Media, Inc., 2009.

[12] M. L. Collard and J. I. Maletic, "Document-oriented source code transformation using XML," in *Proc. SET*, vol. 9, 2004, pp. 11–14.

[13] H. Liu, Q. Liu, C.-A. Staicu, M. Pradel, and Y. Luo, "Nomen est omen: Exploring and exploiting similarities between argument and parameter names," in *Proc. ICSE*, 2016, pp. 1063–1073.

[14] N. Madani, L. Guerrouj, M. D. Penta, Y.-G. Gueheneuc, and G. Antoniol, "Recognizing words from source code identifiers using speech recognition techniques," in *Proc. CSMR*, 2010, pp. 68–77.

[15] B. Caprile and P. Tonella, "Restructuring program identifier names," in *Proc. ICSM*, 2000, pp. 97–107.

[16] D. Lawrie, D. Binkley, and C. Morrell, "Normalizing source code vocabulary," in *Proc. WCRE*, 2010, pp. 3–12.

[17] F. Deissenboeck and M. Pizka, "Concise and consistent naming," *Software Quality Journal*, vol. 14, no. 3, pp. 261–282, 2006.

[18] A. Peruma, M. W. Mkaouer, M. J. Decker, and C. D. Newman, "Contextualizing rename decisions using refactorings, commit messages, and data types," *Journal of Systems and Software*, vol. 169, no. 110704, pp. 1–22, 2020.

[19] V. Arnaoudova, L. M. Eshkevari, M. Di Penta, R. Oliveto, G. Antoniol, and Y. Guéhéneuc, "REPENT: Analyzing the nature of identifier renamings," *IEEE Transactions on Software Engineering*, vol. 40, no. 5, pp. 502–532, 2014.

[20] A. Peruma, M. W. Mkaouer, M. J. Decker, and C. D. Newman, "An empirical investigation of how and why developers rename identifiers," in *Proc. IWOR*, 2018, p. 26–33.

[21] C. D. Newman, R. S. AlSuhaibani, M. J. Decker, A. Peruma, D. Kaushik, M. W. Mkaouer, and E. Hill, "On the generation, structure, and semantics of grammar patterns in source code identifiers," *Journal of Systems and Software*, vol. 170, no. 110740, pp. 1–21, 2020.

[22] C. D. Newman, A. Preuma, and R. AlSuhaibani, "Modeling the relationship between identifier name and behavior," in *Proc. ICSME*, 2019, pp. 376–378.

[23] E. Hill, Z. P. Fry, H. Boyd, G. Sridhara, Y. Novikova, L. Pollock, and K. Vijay-Shanker, "AMAP: Automatically mining abbreviation expansions in programs to enhance software maintenance tools," in *Proc. MSR*, 2008, pp. 79–88.

[24] Y. Jiang, H. Liu, J. Zhu, and L. Zhang, "Automatic and accurate expansion of abbreviations in parameters," *IEEE Transactions on Software Engineering*, vol. 46, no. 7, pp. 732–747, 2020.

[25] A. Alatawi, W. Xu, and J. Yan, "The expansion of source code abbreviations using a language model," in *Proc. COMPSAC*, vol. 2, 2018, pp. 370–375.

[26] L. Surafel, A. Lemma, and T. Paolo, "Automated identifier completion replacement," in *Proc. CSMR*, 2013, pp. 263–272.

[27] Y. Kashiwabara, Y. Onizuka, T. Ishio, Y. Hayase, T. Yamamoto, and K. Inoue, "Recommending verbs for rename method using association rule mining," in *Proc. CSMR-WCRE*, 2014, pp. 323–327.

[28] M. Fowler, *Refactoring: Improving the Design of Existing Code.* Addison Wesley, 1999.

[29] E. Murphy-Hill, C. Parnin, and A. Black, "How we refactor, and how we know it," *IEEE Transactions on Software Engineering*, vol. 38, no. 1, pp. 5–18, 2012.

[30] T. Saika, E. Choi, N. Yoshida, A. Goto, S. Haruna, and K. Inoue, "What kinds of refactorings are co-occurred? An analysis of Eclipse usage datasets," in *Proc. IWESEP*, 2014, pp. 31–36.

[31] K. Maruyama and S. Hayashi, "A tool supporting postponable refactoring," in *Proc. ICSE*, 2017, pp. 133–135.

[32] Y. Osumi, N. Umekawa, H. Komata, and S. Hayashi, "Appendix of empirical study of co-renamed identifiers," 2022. [Online]. Available: https://zenodo.org/record/7214226