



Evaluating Performance of Clone Detection Tools in Detecting Cloned Co-change Candidates

Md Nadim, Manishankar Mondal, Chanchal K. Roy

Department of Computer Science, University of Saskatchewan, Saskatoon, Canada

Abstract

Code reuse by copying and pasting from one place to another place in a code-base is a very common scenario in software development which is also one of the most typical reasons for introducing code clones. There is a huge availability of tools to detect such cloned fragments and a lot of studies have already been done for efficient clone detection. There are also several studies for evaluating those tools considering their clone detection effectiveness. Unfortunately, we find no study which compares different clone detection tools in the perspective of detecting cloned co-change candidates during software evolution. Detecting cloned co-change candidates is essential for clone tracking. In this study, we wanted to explore this dimension of code clone research. We used 12 different configurations of nine promising clone detection tools to identify cloned co-change candidates from eight C and Java-based subject systems and evaluated the performance of those clone detection tools in detecting the cloned co-change fragments. Our findings show that a good clone detector may not perform well in detecting cloned co-change candidates. The amount of unique lines covered by clone fragments, the number of detected clone fragments, source file processing mechanism, types of detected clones plays an important role in detecting cloned co-change candidates.

Keywords: Clone Detection, Cloned Co-change Candidates, Commit operation, Software Maintenance.

Email address: {mdn769, mshankar.mondal, chanchal.roy}@usask.ca (Md Nadim, Manishankar Mondal, Chanchal K. Roy)

1. Introduction

A large number of software tools have already been introduced for detecting cloned code fragments. Two surveys, that were done in 2009 by Roy et al. [1] and in 2013 by Rattan et al. [2] reported 75% increase in the number of clone
5 detection tools in these four years. Roy and Cordy [3] reported the existence of about 200 tools for detecting cloned code fragments. Although a large number of clone detection tools currently exist, we found no study for comparing the performance of different tools based on their ability to be used in software main-
10 tenance activity such as predicting cloned co-change candidates during software evolution. In this study, we wanted to explore, whether a good clone detector also performs well in detecting cloned co-change fragments?

One of the common features of clone detection tools is to combine similar code fragments into a clone group or class. The code fragments in a particular clone class are expected to perform similar functionalities. If we want to make
15 changes to a particular clone fragment in a clone class, the other fragments in the class are likely to have similar changes to ensure consistency of the code-base. Considering this assumption, we can say that all the clone fragments in a clone class have the possibility of being a cloned co-change candidate with any change of that class members. We utilize the clone classes provided by the clone
20 detectors for these types of co-change prediction.

Finding the co-change candidates of a target code fragment is also known as change impact analysis [4] in the literature. Mondal et al. [5] investigated whether a clone detection tool can enhance the performance of an evolutionary coupling based tool in finding change impact set or co-change candidates. They
25 performed their investigation using Nicad for detecting both the regular and micro-clones and found that use of detected clone results significantly enhance the performance of Tarmaq Rolfsnes et al. [6]. As they only analysed the use of Nicad, in this study we wanted to compare some other good clone detection tools to find whether these tools can perform better for detecting co-change

30 candidates. We have evaluated four different configurations of CloneWorks [7] and eight other clone detectors in our investigation. Therefore, we have a total of 12 separate implementations of clone detection tools (we will consider them as 12 separate tools in the rest of this paper). We apply these tools on eight open-source software systems. Configuration of the clone detection tools are 35 given in Table 3 and the software systems used in this study are reported in Table 1.

During software evolution, a developer makes changes in the code-base to fulfil some change requests. Those change requests could be related to each other or independent [8, 9]. Therefore, all the changes done in a single commit 40 need not be related to each other. Some changes in a single commit may be dependent on each other and some may be independent. The related code fragments are known as the co-change candidates in literature [10]. Some of those co-change candidates may contain similar code-fragments i.e. they are clones of one another, on the other than, other types of co-change candidates 45 may not be cloned fragments but they have a functional dependency or coupling with each other. If a developer makes changes to a target code fragment, those changes might also need to be reflected other similar fragments in the code-base to ensure consistent evolution of the software system [5, 11]. Failing to change a co-change candidate of a target fragment can introduce bugs in the software 50 system [12, 13].

We have analyzed thousands of commit operations from the evolutionary histories of eight subject systems listed in Table 1. While analyzing a commit operation, we identify which code fragments changed together (i.e., co-changed) in that commit. Considering each fragment as the target fragment, we try to 55 predict the other actually co-changed fragments using each of our clone detectors. We found some change fragment which is not detected by any of the clone detectors. We excluded those change fragments from consideration during calculating the performance measures of clone detectors. An example of our detection process is demonstrated in Fig. 1. Let us assume that 21 changes, C1

60 to C21, occurred in the code-base of a subject system in a particular commit operation. We detect these changes using the UNIX diff operation. If we consider C1 as the target change, the other 20 changes, C2 to C21, are the actual co-change candidates (i.e., co-changed candidates) of C1. We apply different clone detectors to detect these co-change candidates for the target change C1.

65 Let using Deckard we can detect five change fragments (C2, C6, C8, C15, C21) from those 20 fragments, similarly using Nicad we can detect four fragments (C5, C10, C16, C18). We will continue to detect co-change fragments using all the other clone detectors. After getting the results from all the clone detectors, we find 10 unique change fragments (C2, C5, C8, C15, C21, C5, C8, C10,

70 C15, C21) out of 20 fragments by taking a union of the results of all the clone detectors. We will take those 10 unique change fragments as cloned co-change candidates and calculate the precision and recall of each of the clone detectors based on their number of detection among those cloned co-change candidates. For each subject system, we finally calculated average recall, average precision,

75 and F1 Score for each of the clone detector and then compare the clone detectors based on their weighted average F1 Score considering all the subject systems in this study. Figure 2 shows the bar chart of Average Recall and Average Precision drawn from our experimental results and in Table 5 we have given the calculated weighted average of F1 Score. According to our findings and ranking

80 of the clone detectors (Table 6 and 7), we can conclude that CloneWorks (both two configurations, Type-3 Pattern and Token), Deckard, and CCFinder outperforms all the other tools. CloneWorks Type-2 Blind, ConQAT, and iClones falls in the following order. From the final rank list, we also see that the clone detection tools which detecting only Type-1 clones (such as Duplo, CloneWorks

85 Type-1) are performing worst in finding co-change candidates. We also calculated the average number of distinct lines detected as cloned lines by each of the clone detectors in all the revisions of all the subject systems (Figure 3) and found that the clone detector which detects more distinct lines as a cloned line in the code-base also performs well in detecting cloned co-change candidates.

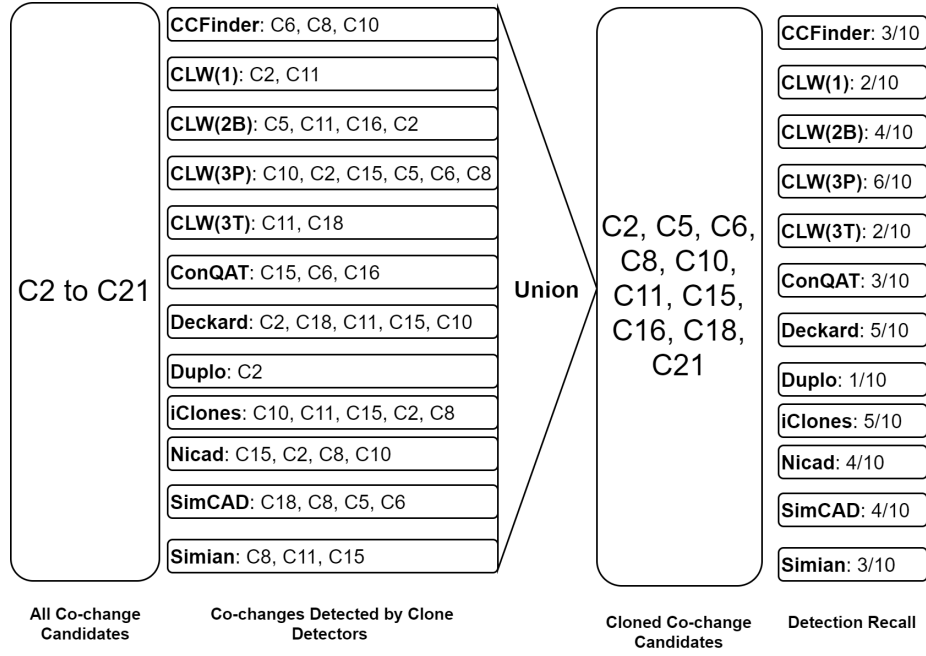


Figure 1: Demonstrating cloned co-change detection process

Based on this study, we tried to answer the following research questions:

RQ1: What is the comparison scenario of the clone detectors in predicting cloned co-change candidates?

RQ2: Why do different clone detectors perform differently in detecting cloned co-change candidates?

RQ3: Do the source code processing techniques (Pattern/Token/Text based processing) of the clone detection tools have any impact on their performance in detecting co-change candidates?

RQ4: Do clone detection tools designed for detecting different types of clones (Type 1, 2, 3) work differently in detecting cloned co-change candidates?

To the best of our knowledge, our study is the first one to compare clone detection tools considering a particular maintenance perspective (e.g., considering their capabilities in successfully suggesting cloned co-change candidates

during software evolution). From an initial assumption, it is obvious that a clone detector which is good in detecting clone should also be good in detecting cloned co-change candidates. In this exploratory study, we wanted to practically verify this assumption. We selected 12 implementations of clone detectors detecting different types of clones in our investigation to verify whether they are also good at detecting co-change candidates. According to our investigation and analysis, we find that, the clone detectors which detects more type-3 clones and performs pattern based source code processing are good in detecting cloned co-change candidates. Our investigation also shows that, tools which provides more number of clone fragments and covers more source code lines are also good in detecting cloned co-change candidates. We have created two rank lists of the clone detection tools based on our investigation. One of them in Table 6 where we have considered the ranking of the clone detectors in each of the subject systems. A clone detector which performed good in most of the subject system got the higher rank in this ranking table. Other ranking is in Table 7 where we considered the weighted average (number of changes whose co-change candidates are detected are the weight factor) of F1 Scores of clone detectors considering all the subject systems. Considering both the rankings we find common tools obtaining in top-5 positions. Three of them are the different configurations of CloneWorks except Type-1, and other two clone detectors are Deckard and CCFinder. Therefore, we can conclude that, to detect cloned co-change candidates, those five tools (all are pattern and token based) are the best choice compared to the other seven tools used in this study. From these ranking, we can also find that, text based clone detectors (such as Duplo or CloneWorks Type-1 configuration) are not good for detecting co-change candidates.

This paper is a significant extension of our previous work [14] on detecting cloned co-change candidates using different clone detectors. While in our previous study, we have used six open-source software as subject systems to evaluate the performance of six clone detection tools which have been reported good in detection clone fragments from software source code by recent studies. By answering two research questions, we have been evaluated that, even though

a tool which is good in detecting clone fragments from software systems may not be good in detecting cloned co-change candidates. We also investigated and reported two possible reasons for such a difference in the performance of clone detection tools while we are using them to predict co-change fragments. We extend our previous work by answering two additional research questions (RQ3, RQ4) to find a more specific reason for the variation of the performance by clone detectors in detecting co-change candidates. We have also increased the generalizability of the previous study by adding two more software systems as subject system totalling the number of subject systems in six. We have added three more clone detectors (totalling nine clone detectors) including four different configurations of one clone detector (CloneWorks) which makes 12 total clone detector executions. We used different configurations of CloneWorks for targeting different types (Type-1, or Type-2 blind, or Type-3 pattern, Type-3 token) of clones to investigate the effect of those types in our study. Therefore, our implementation has been upgraded from 6X6 to 12X8 (Clone detector X Subject Systems) in the current version of the study. In the current extended version of the study, we have shown that the performance of clone detection tools in detecting cloned co-change fragments not only dependent on the number of clone fragments detected and the line covered in the source file by those fragments but also the type of detected clone and underlying source code processing techniques also have some impact.

We organized this paper in the following sections: Some related works are described in Section 2, our methodology is in Section 3, we described the experimental result in Section 4, the discussion is in section 5, Section 6 explains some possible threats to validity, and we conclude our paper by mentioning future work in Section 7.

2. Related Work

There are several studies [15, 1, 16, 17] that have been focused on ranking different clone detection tools based on their performance and accuracy in de-

tecting different types of clone fragments. Burd and Bailey [18] did a study for
 comparing the performance of three clones and two plagiarism detecting tools
 165 based on their precision and recall of the ability to detect duplicated codes in
 a single file or across different files. Bellon et al. [16] evaluated six clone de-
 tection tools based on eight large C and Java programs of almost 850 KLOC
 and made a framework for comparing different clone detection tools with the
 data validated by one of its authors. Rysselberghe and Demeyer [19] evaluated
 170 three representative clone detection techniques from a refactoring perspective
 where they provided comparative results in terms of portability, kinds of clone
 reported, scalability, number of false positive, and number of useless clone de-
 tection. Svajlenko and Roy [15] evaluated eleven modern clone detection tools
 using four benchmark frameworks and noted ConQAT, iClones, NiCad and
 175 SimCAD as very good tools for detecting clones of all the three types (Type-1,
 Type-2, Type-3). Roy et al. [1] did a qualitative comparison and evaluation of
 the latest clone detection approaches and tools, and made a benchmark called
 BigCloneBench [3] which contains eight million manually validated clone pairs
 in a large inter-project source dataset of more than 25,000 projects and 365
 180 million lines of code. They categorize, relate and assess different clone detection
 tools based on two different points of view such as classification based on the
 overlapping set of attributes in the different code fragments and the scenarios
 how Type-1, Type-2, Type-3, and Type-4 clones created. They also elaborated
 the procedure of using the result of their study to select the most suitable clone
 185 detection tool or technique in the context of a specific set of areas and limita-
 tions.

There are some studies which not only proposed a clone detection mechanism
 but also did a comparison of their proposed technique with some existing tech-
 niques. Koschke et al. [20] provided a technique to detect clone using suffix
 190 trees in abstract syntax trees and they also made a comparison to other tech-
 niques using the Bellon benchmark for clone detectors. Ducasse et al. [21] and
 Selim et al. [22] also utilized Bellon's framework for measuring the performance
 of their proposed clone detection tools based on string comparison and inter-

mediate source transformation respectively. Selim et al. [22] showed that their
195 tool is capable of detecting Type-3 clones and their technique is better than the
source-based clone detectors based on the value of recall through a slight drop
in the precision using Bellon’s corpus where clone group is not complete. Com-
pared to the standalone string and token-based clone detectors, their technique
showed a little higher precision.

200 All the studies which compared different clone detectors have been focused on
the precision, recall, computational complexity, and memory used or detecting
a specific type of clone fragments such as Type-1, Type-2, Type-3, or Type-4
during the detection approach of duplicated code in a code-base. Our study to
compare clone detectors is completely different from the previous comparisons.

205 We do not want to compare clone detection tools based on the capability to
detect clones. Our point of interest is to detect co-change candidates during the
software commit operations. Mondal et al. [10] did a study to predict and rank
the co-change candidates by analyzing evolutionary coupling from previously
done change history using generated clone fragments by NiCad but they did
210 not consider the result of other clone detection tools and also did not show any
comparative study among different clone detection tools in doing such predic-
tion and rank of co-change candidates. This work is an extended version of our
previous study [14] using six clone detection tools on six software systems writ-
ten in C and Java programming languages to compare those tools based on the
215 performance of detecting clone co-change candidates. We found no other study
which have performed similar comparison of clone detectors. To extend our pre-
vious research, we have analyzed the performance of nine clone detection tools
in 12 different configurations based on their capabilities in finding co-change
candidates during software evolution using their generated clone results. Ac-
220 cording to our knowledge, this is the first such investigation of performance with
clone detection tools.

3. Methodology

We have used eight subject systems listed in Table 1 and 12 clone detection tools (Table 3) for our analysis. Our analysis aims to rank these clone
225 detection tools based on their performance in successfully suggesting actual co-change candidates (ACC) during the software evolution. Before starting our main analysis, we have to resolve some issues and we have taken the following considerations in this regard.

Selection of subject systems: To select subject systems for this study, we
230 considered both the popularity of programming language and availability of a considerable amount of revisions. According to the TIOBE Programming Community index [23] (an indicator of the popularity of programming languages), Java is dominating the list of popular programming languages for more than the last ten years and C is the second most popular programming language within
235 this period. Considering this fact, we wanted to select subject systems written in these two programming languages. Our other consideration was the availability of a considerable amount of revisions of each of the systems. Based on both of the considerations, we have chosen the subject systems listed in Table 1. Four of our eight subject systems are written in C programming language
240 and other four are in Java. To increase the generalizability of the study we have added systems having diverse size and application domains.

Selection of clone detectors: In this research, we wanted to examine those clone detection tools which are good in detecting all types of clones. To select such tools, we considered some related studies. We have taken CloneWorks
245 [7] as it is considered as a fast and flexible clone detector for large-scale near-miss clone detection experiments. CloneWorks tool provides the ability to change its processing mechanism by changing its configuration files. We applied four different configurations of CloneWorks to detect Type-3 Pattern, Type-3 Token, Type-2 Blind, and Type-1 clones for investigating the impact of the types
250 of clones in detecting co-change candidates. We included Duplo [24] as another type-1 clone detector for making the comparison with type-1 clones of

Table 1: SUBJECT SYSTEMS

Systems	Language	Domains	Revisions
Brlcad	C	Computer Aided Design	2115
Camellia	C	Batch Job Server	301
Carol	Java	Game	1700
Ctags	C	Code Def. Generator	774
Freecol	Java	Game	1950
Jabref	Java	Reference Manager	1545
jEdit	Java	Text Editor	4000
Qmailadmin (QMA)	C	Mail System Manager	317

Table 2: SUMMARY OF DATA PROCESSED

Revisions/ Subject Systems	Brlcad	Camellia	Carol	Ctags	Freecol	Jabref	jEdit	QMA
Processed	2113	301	1700	774	1001	1540	215	317
Experiencing change	660	163	454	447	836	860	145	35
Experiencing more than one change	553	155	430	330	833	755	145	25

CloneWorks in this investigation. ConQAT [25], iClones [26], NiCad [27], and SimCAD [28] have been reported as very good tools for detecting all types of clones in the study of Svajlenko and Roy [15]. Besides these, CCFinder [29],
 255 Deckard [30], iClones and NiCad are often considered as common examples of modern clone detectors that support Type-3 clone detection. CCFinder is known as a multi-linguistic token-based code clone detection system for large scale source code. Inclusion of CCFinder enriched the variation of detected clone fragments in the extended study. To make more comparison of the per-
 260 formance of type-1 clones in detecting co-change candidates we added Duplo in our study. The reason of taking Simian [31] in our analysis was its ability to find duplicated code by line-by-line textual comparison supporting identifier renaming with a fast detection speed on the large repository and extensive use

in several clone studies [32, 33, 34, 35, 36]. NiCad, SimCAD, and Simian are
265 textual similarity-based clone detection tools. Deckard works using tree compar-
ison based technique. CCFinder, ConQAT and iClones are token-based clone
detection tools.

**Determining if the extracted co-changes are related to each other
or not:** Even though we have extracted all the changes between two adjacent
270 revisions (i.e., revision n and $n+1$), it is not possible to fully guarantee that all
the changes are actually co-change candidates of each other. There might be
some changes which do not depend on any other changes i.e. they may change
independently. The inclusion of such dissimilar changes into our calculation
can drop the detection accuracy of clone detectors. To minimize such drops,
275 we excluded those co-changes which are not detected by any of the 12 clone
detection techniques in our study. As none of the clone detectors considers
them as co-change candidates, we considered those changes as dissimilar or
independent changes.

**Ensuring if the configuration parameters of all the clone detection
280 tools identical with each other or not:** As we wanted to compare different
clone detectors based on their capability of successfully suggesting co-change
candidates, it was important to configure them identically during detecting
clones from our subject systems. Wang et al. [33] introduced confounding con-
figuration choice problem where the configuration of different tools during clone
285 detection may play a vital role and the result may be best or worst depending
on the configuration. Our configuration of different tools is shown in Table 3.
We have used similar configurations for each of the tools for obtaining a consis-
tent result. We have taken configuration values similar to Svajlenko and Roy
[15] which they conducted to compare different clone detectors based on their
290 efficiency in detecting cloned fragments. We provided 70% similarity threshold
for all the clone detection tools (except Deckard) which takes similarity dissim-
ilarity value as parameter. We have used 85% as the similarity threshold for
Deckard 85% because we found a lot of unwanted clones in the result if we use

the similarity threshold 70%. These results include a lot of duplicated clone
 295 fragments and showing a lot of fragments as a clone to itself several times. We
 also tried some other percentage values such as 75%, and 80% but the detected
 result of Deckard becomes much desirable when we set it to 85%. Svajlenko and
 Roy [15] also used 85% similarity while running Deckard for Mutation Frame-
 work. We have also selected identical parameter values such as the minimum
 300 number of tokens, minimum number of lines for different clone detection tools.
 As we wanted to compare different clone detectors based on their capability of
 successfully suggesting co-change candidates, it was very important to configure
 them identically during detecting clones from our subject systems.

The overall approach: Our overall processing is performed in some dis-
 305 tinct steps. Initially, we downloaded all the source files of all the revisions of
 all the subject systems from their respective SVN repositories. We then ap-
 plied **diff** operation between each file of a revision with the respective file in the
 next revision and extracted the change information such as Name of the File
 which is changed, the Line where the respective change begins, the Line where
 310 the change is ended from the output of **diff**. We did the change extraction for
 each of the revision (excluding the last one) of all the subject systems. After
 detecting all the changes, we started the clone detection on all the revisions of
 all the subject systems using all the clone detection tools. We started our main
 analysis to find the accuracy of each of the clone detection tools after having the
 315 result of all the clone detectors and change information from all the revisions.

The mechanism of calculating accuracy is demonstrated in our introduction
 using Fig. 1. Suppose, we are examining a particular commit operation. The
 number of fragments that were changed in this commit operation is n . Now,
 let us consider one of these n fragments as the target fragment. Then the other
 320 $n - 1$ fragments are the actually co-changed candidates for the target fragment.
 We excluded the non-cloned co-change candidates using the approach described
 in the introduction. After this exclusion, we get the **Actually Cloned Co-
 change** (ACC) for each of the target fragments.

Let us assume that the target change fragment intersects a particular clone

325 fragment from a particular clone class. The other fragments in that clone class
 are considered as the **Predicted Cloned Co-change** (PCC) candidates. We
 now determine how many of these PCC intersect with the ACC to obtain the
 number of detected cloned co-change candidates by the clone detector.

These counts of predicted and actually co-changed candidates are considered
 330 as the **true positives** to calculate Recall, Precision, and F1 Score. We calculate
 these using the following equations (Eq. 1, 2, and 3).

$$Recall = \frac{|PCC \cap ACC|}{|ACC|} \quad (1)$$

$$Precision = \frac{|PCC \cap ACC|}{|PCC|} \quad (2)$$

$$F1 \text{ Score} = \frac{2 \times Precision \times Recall}{Precision + Recall} \quad (3)$$

We repeat the calculating process of Recall and Precision for all the changes
 in each of the subjects systems with the detected clone fragments generated by
 all the clone detection tools. We then calculate F1 Score of the clone detectors
 335 for each of the subject systems by taking the average values of Recall and
 Precision which is reported in Table 5. We reported both the ranking of the
 tools considering ranks in each of the subject systems in Table 6 and ranking
 considering weighted average of F1 Score in all the subject systems in Table 7.
 To calculate the weighted average of F1 Score of the tools we took the total
 340 number of changes is the corresponding weighting factor in each subject system
 whose co-change is detected by any one of the clone detectors in this study.

4. Experimental Result

In this section, we will answer the research questions based on our overall
 analysis and obtained results by processing each of the eight subject systems
 345 using all the 12 clone detection tool executions.

Table 3: CONFIGURATION OF PARTICIPATING CLONE DETECTION TOOLS

Tools	Configuration for Clone Detection
CCFinder	min. size: 50 tokens, min. token types: 12
CLW(T1)	termsplit=token, termproc=Joiner
CLW(T2Blind)	cfproc=rename-blind, cfproc=abstract literal, termsplit=token, termproc=Joiner
CLW(T3Pattern)	cfproc=rename-blind, cfproc=abstract literal, termsplit=line
CLW(T3Token)	termsplit=token, termproc=FilterOperators, termproc=FilterSeperators
ConQAT	block clones, clone min-length=5, gap ratio=0.3
Deckard	min. size: 30 tokens, 5 token stride, min. 85% similarity
Duplo	min. size: 10 lines, min. characters/line:1
iClones	minimum block: 30, minimum clone: 50, All Transformation
Nicad	block clones, blind renaming, max. threshold=0.3, minimum lines=5, maximum lines=2500
SimCAD	block clones, Source Transformation= generous
Simian	min. size: 5 lines, normalize literals/identifiers

CLW: *CloneWorks*, **T1/2/3:** *Type-1/2/3*

Table 4: SUMMARY OF ACTUAL TARGET AND CO-CHANGE CANDIDATES

SS	# ATC	# ACC	% ATC	% ACC
Brlcad	2909	33578	7.45	1.89
Camellia	8052	346140	20.61	19.46
Carol	4582	254311	11.73	14.29
Ctags	718	3648	1.84	0.21
Freecol	6865	265213	17.57	14.91
Jabref	8313	455469	21.28	25.60
jEdit	5122	323277	13.11	18.17
QMA	2508	97396	6.42	5.47
Total	39069	1779032	100	100

* *SS: Subject Systems*

* *# ATC: Number of Actual Target Changes*

* *# ACC: Number of Actual Co-changes*

4.1. Answer to the *RQ1*

What is the comparison scenario of the clone detectors in predicting cloned co-change candidates?

350 The key experimental results are in Figure 2, Table 4, Table 5, Table 6, and Table 7 where Fig. 2 shows the average Recall and average Precision of each of the clone detection tools. Table 4 shows the summary of target changes and detected co-change candidates for those target changes in each of the subject systems. We found the highest and lowest percentage of target change and its
355 cloned co-change candidates from Jabref and Ctags respectively. Table 5 shows the F1 Score of each of the clone detectors in each of the subject systems. The F1 Score is calculated using Equation (3). Our experimental results concluded in Table 6 which shows that CLW(T3Pattern), CLW(T3Token), and Deckard shows top performance (Rank 1 or 2) in most of the subject systems compared to
360 all the other tools. The summary of the results in the Table 6 shows that among

the subject systems, CLW(T3Pattern) is the best in all the subject systems except Camellia and Freecol where Deckard is showing the best performance. CLW(T3Token) shows the second-best performance in most of the subject systems. An overall observation on individual rankings of different clone detection techniques reveals that CLW(3Pattern), Deckard, CLW(3Token), CCFinder show better performance in most of the subject systems compared to the other clone detectors. On the other hand, Duplo, CLW(Type1) shows worst performance in most of the subject systems. Other tools shows average performance considering individual ranking in different subject systems. We also calculated overall ranking considering weighted average of F1 Scores in individual subject systems which is reported in Table 7. Though, in individual ranking in Table 6 CLW(3Pattern) showed best performance compared to Deckard, weighted average of Deckard become higher in Table 7 as Deckard obtained much better F1 Score in Camellia than CLW(3Pattern). Camellia has a very high number of target changes (#ATC in Table 4) which makes weighted average of F1 Score higher. CLW(3Pattern) obtained second-highest position in the rank list by obtaining very close weighed average F1 Score to Deckard. CLW(Type1) and Duplo obtained the bottom position in the final rank list.

As our analysis was based on the clone grouping into class or pair provided by the clone detection tools, we found that the efficiency of clone detection tools in suggesting cloned co-change candidates is mostly dependent on its effectiveness in making clone class/ pair. The tool which groups functionally similar clone fragments into a clone class/ pair effectively can perform well in successfully suggesting cloned co-change candidate(s). Different values of the accuracy of different clone detectors indicate the difference in their efficiency in this research domain.

4.2. Answer to the **RQ2**

Why do different clone detectors perform differently in detecting cloned co-change candidates?

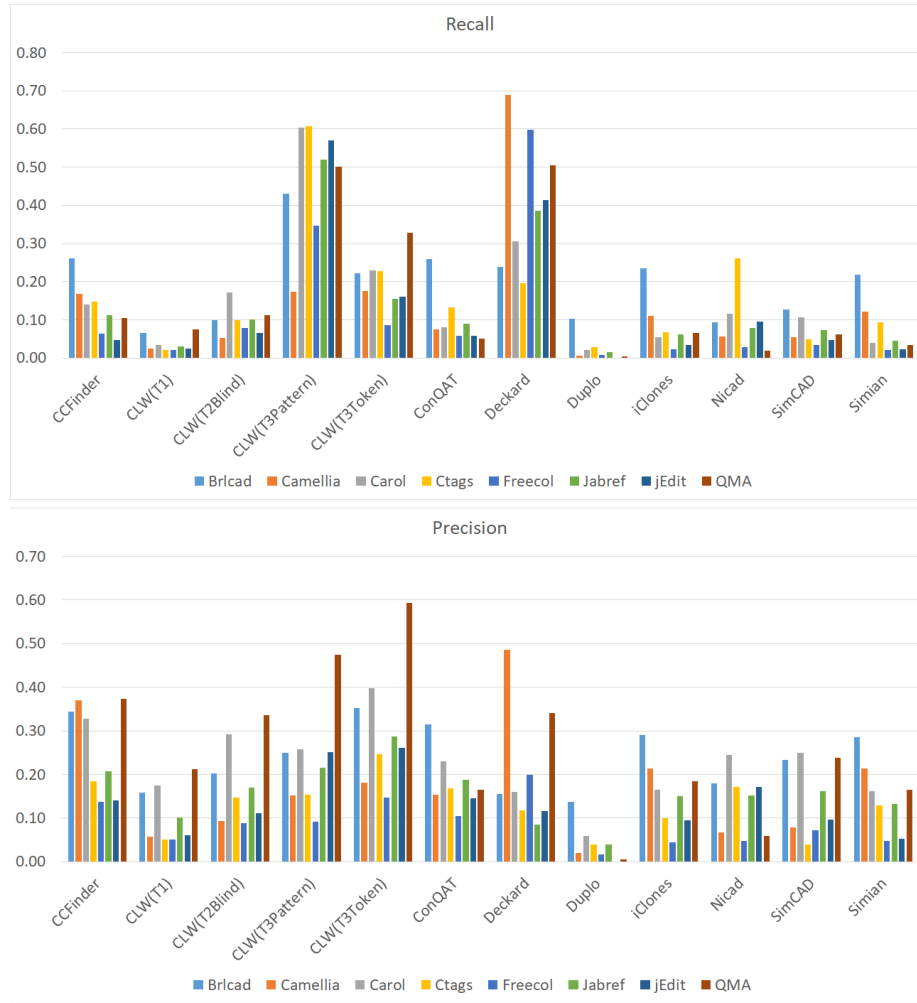


Figure 2: Average recall of different tools

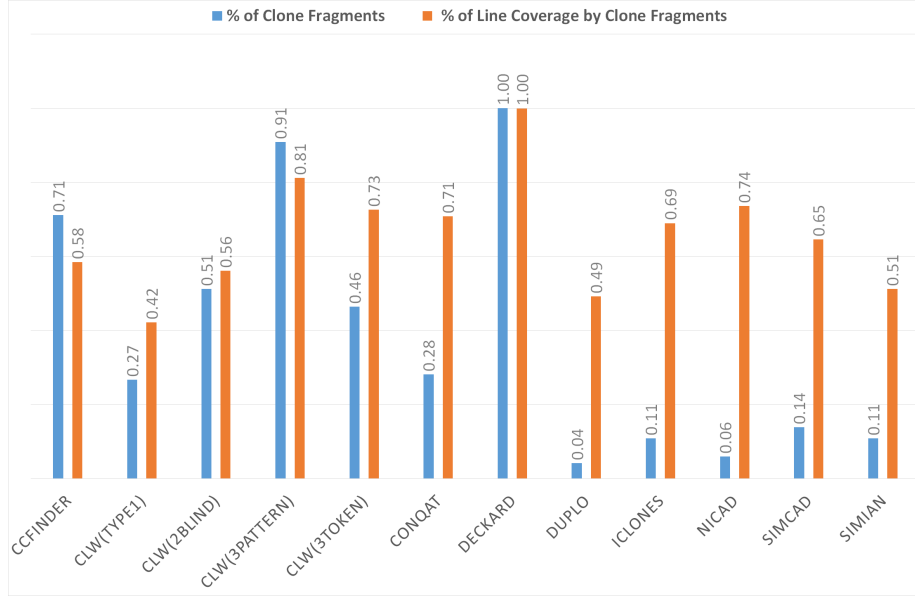


Figure 3: Comparing unique line coverage by clone fragments and number of clone fragments from different clone detectors.

From the answer of our **RQ1**, we found a difference in performance for different clone detection tools in suggesting cloned co-change candidates. We found a clone detection tool which is good in detecting clone fragments may not be good at detecting cloned co-change candidates. This motivates us to find out the reason to answer this research question.

We investigated the number of clone fragments and the number of unique lines covered by those clone fragments by all the 12 clone detectors from all the revisions of all the subject systems. Figure 3 shows the comparison scenario of number of clone fragments and line covered by those clone fragments from different clone detectors. For better comparison, we bring the values in a single scale (between 0 and 1) where 0 and 1 represents the lowest and highest values respectively compared to all the clone detectors under comparison. Considering both, the the number of clone fragments and the number of lines covered by those clone fragments from all the revisions of all the subject systems, if we order the clone detectors from the highest to the lowest, we find Deckard and

CLW(3Pattern) in the top of the list. CLW(3Token) and CCFinder fall in the respective next position in providing highest number of clone fragments and covering highest number of unique lines in the source files. This scenario shows that, a good clone detector can perform bad in detecting cloned co-change candidates if it does not detect enough clone fragments and does not cover enough unique lines by those clone fragments in the source file. Though, earlier study [10] suggests that NiCad is a very good clone detector, in both of these cases, it falls at the bottom of the list. Despite, NiCad performs very well in detecting clone fragments, it provides a lower number of clone fragments and also the lower number of line coverage by those clone fragments in the software systems. For that reason, while detecting the cloned co-change candidates, NiCad is showing lower F1 Score. The number of clone fragments and line coverage by those fragments seems to be an underlying factor behind the obtained comparison scenario of the clone detectors in predicting cloned co-change candidates, there can be several other factors such as overlapping of code clones and code similarity detection mechanism. We plan to investigate these factors in future.

4.3. Answer to the **RQ3**

Do the source code processing techniques (Pattern/Token/Text based processing) of the clone detection tools have any impact on their performance in detecting co-change candidates?

We can answer this research question by analysing our final ranking of the clone detectors in Table 6 and 7. Top two clone detectors (Rank 1 and 2) work by extracting source code patterns from the code-base. Deckard first generates vectors from the source file and then extracts tree like source pattern to match similarity among different source code fragments. CLW(3Pattern) also processes the source code terms by splitting in lines and then extracts code patterns. The other five tools (Rank 3 to 7) in the rank list perform token-based source code processing and the remaining five tools perform text based source code processing for detecting clones from the source file. From this result we can say that, text based clone detection tools are not good to be used in detecting

Table 5: F1 SCORE OF DIFFERENT TOOLS IN DETECTING CLONED CO-CHANGE

Tools/SS	Brlcad	Camellia	Carol	Ctags	Freecol	Jabref	jEdit	QMA
CCFinder	0.30	0.23	0.20	0.16	0.09	0.15	0.07	0.16
CLW(T1)	0.09	0.04	0.06	0.03	0.03	0.05	0.04	0.11
CLW(T2Blind)	0.13	0.07	0.22	0.12	0.08	0.13	0.08	0.17
CLW(T3Pattern)	0.32	0.16	0.36	0.25	0.15	0.30	0.35	0.49
CLW(T3Token)	0.27	0.18	0.29	0.24	0.11	0.20	0.20	0.42
ConQAT	0.28	0.10	0.12	0.15	0.08	0.12	0.08	0.08
Deckard	0.19	0.57	0.21	0.15	0.30	0.14	0.18	0.41
Duplo	0.12	0.01	0.03	0.03	0.01	0.02	0.00	0.00
iClones	0.26	0.15	0.08	0.08	0.03	0.09	0.05	0.10
Nicad	0.12	0.06	0.16	0.21	0.04	0.10	0.12	0.03
SimCAD	0.17	0.07	0.15	0.04	0.05	0.10	0.06	0.10
Simian	0.25	0.16	0.06	0.11	0.03	0.07	0.03	0.06

cloned co-change candidates during software evolution. The tools which can detect more generalized clone fragments specially pattern based clone detectors are very good for detecting co-change candidates.

4.4. Answer to the **RQ4**

440 **Do clone detection tools designed for detecting different types of clones (Type 1, 2, 3) work differently in detecting cloned co-change candidates?**

From the final rank list of our clone detectors we also find the relation of detected clone types with its ability to detect cloned co-change candidates.
 445 Rank list of clone detectors in Table 6 and 7 shows that, clone detecting tools such as CLW(Type1), Duplo, which detects only Type 1 clone will not perform good in detecting co-chagne candidates. On the other hand, tools such as Deckard, CLW(3Pattern), CLW(3Token), CCFinder perform very good in detecting cloned co-change candidates.

Table 6: RANKS OF THE CLONE DETECTORS BY CONSIDERING INDIVIDUAL RANKING IN EACH OF THE SUBJECT SYSTEMS

Tools/ Subject Systems	S1	S2	S3	S4	S5	S6	S7	S8	Sum of Ranks	Final Rank
CLW(T3Pattern)	1	4	1	1	2	1	1	1	12	1
CLW(T3Token)	4	3	2	2	3	2	2	2	20	2
Deckard	7	1	4	6	1	4	3	3	29	3
CCFinder	2	2	5	4	4	3	7	5	32	4
CLW(T2Blind)	9	8	3	7	5	5	5	4	46	5
ConQAT	3	7	8	5	6	6	6	9	50	6
iClones	11	10	6	3	8	7	4	11	60	7
Simian	5	6	9	9	10	9	9	7	64	8
Nicad	8	9	7	10	7	8	8	8	65	9
SimCAD	6	5	11	8	11	10	11	10	72	10
CLW(T1)	12	11	10	11	9	11	10	6	80	11
Duplo	10	12	12	12	12	12	12	12	94	12

** S1-S8 represents sequence of eight subject systems used in this study.*

Table 7: RANKS OF THE CLONE DETECTORS BY CONSIDERING WEIGHTED AVERAGE F1 SCORE
IN ALL THE SUBJECT SYSTEMS

Clone Detectors	WA of Precision (p)	WA of Recall (r)	F1 Score $2pr/(p+r)$	Final Rank
Deckard	0.22	0.47	0.30	1
CLW(3Pattern)	0.21	0.43	0.28	2
CLW(3Token)	0.27	0.17	0.21	3
CCFinder	0.25	0.12	0.16	4
CLW(2Blind)	0.16	0.09	0.12	5
ConQAT	0.17	0.09	0.12	6
iClones	0.15	0.07	0.10	7
Nicad	0.13	0.07	0.09	8
Simian	0.14	0.07	0.09	9
SimCAD	0.14	0.07	0.09	10
CLW(Type1)	0.10	0.03	0.05	11
Duplo	0.03	0.02	0.02	12

450 5. Discussion

There are two primary perspectives of managing code clones: (1) clone tracking and (2) clone refactoring. Our research essentially focuses on the clone tracking perspective. The main task of a clone tracker is to suggest similar co-change candidates when a programmer attempts to change a code fragment. For suggesting co-change candidates, a clone tracker depends on a clone detector. Our research compares 12 promising clone detectors based on their capabilities in suggesting cloned co-change candidates. According to our investigation, CloneWorks (Type-2 and 3), Deckard, and CCFinder are the most promising tools for suggesting such co-change candidates based on the both rankings we obtained in Table 6 and Table 7. Based on our overall observation, we can say that the performance of CloneWorks (Type-3 Pattern/ Token) and Deckard are much better compared to the other clone detection tools in detecting co-change candidates during software evolution. As the clone classes/ pairs generated by different clone detectors played an important role in our analysis, we can say that the clone detectors which can group similar clone fragments into a class efficiently will perform better in detecting co-change candidates during the commit operation. From our findings we can also say that the clone detectors which detect all the Type 1, 2, 3 clones are also perform good in detecting co-change candidates.

470 When a particular code fragment is changed, we apply the clone detectors to predict which other similar code fragments might also need to be co-changed. However, some dissimilar fragments might also be changed together with the particular fragment. As we are applying only clone detectors, we cannot consider those dissimilar co-change candidates in our research.

475 In our research, we do not compare the clone detectors considering their clone detection efficiency. We rather compare the clone detection tools based on their ability in suggesting cloned co-change candidates. Such a comparison of clone detectors focusing on a particular maintenance perspective was not done previously. Suggesting co-change candidates for a target program entity is an

480 important impact analysis [4] task during software evolution. Thus, through
our research, we investigate which of the clone detectors can be useful in change
impact analysis to what extent. Findings from our research can identify which
clone detector(s) can be promising for change impact analysis.

6. Threats to Validity

485 We have investigated eight subject systems in our study. While more sub-
ject systems could generalize our findings, we selected our systems focusing on
their diversity, popularity of used programming language, and availability of
a considerable number of revisions. For example, our systems are of different
application domains, sizes, and revision history lengths. Thus, our findings are
490 not biased by our choice of subject systems. We believe that our findings are
important from the perspectives of software maintenance.

We have investigated 12 different configuration of nine clone detectors in our
study. Detection parameter settings of the clone detectors can have an impact
on their comparison. However, the parameters of different clone detectors were
495 selected considering their equivalence. Thus, we believe that we have a fair
comparison among the clone detectors.

Several code fragments might change together in a commit operation. While
some of these fragments can be similar to one another, and some might be
dissimilar. Similar code fragments co-change (i.e., change together) for ensuring
500 consistency of the code-base. However, dissimilar code fragments can co-change
because of their underlying dependencies which could have some impact on
the generalization of this research outcome. As we aim to compare the clone
detection tools, we wanted to discard the dissimilar co-change candidates from
our consideration. If a co-change candidate was not detected as a true positive
505 by any of the clone detectors, we discarded the candidate. We believe that such
a consideration is reasonable in our experiment aiming towards comparing clone
detectors and our findings may inspire more similar research.

7. Conclusion and Future Works

In this research, we make a comparison among different clone detection tools from the perspective of software maintenance. In particular, we investigate their performances in successfully suggesting (i.e., predicting) cloned co-change candidates during evolution. We used eight open source subject systems written in C and Java for our analysis. According to our findings (Table 6 & 7) on thousands of revisions of these systems, CloneWorks (Type-3) and Deckard show the most promising results in most of the eight subject systems compared to the other tools. CCFinder and CloneWorks (Type-2) are also shows better performance in some of the subject systems. Although we have figured some reasons of the better performance of Deckard, CloneWorks, and CCFinder in the Discussion section of our study, we planned to extend this research by analyzing the clone detection mechanism of the clone detectors to find out some other reasons for their performance. We also want to investigate the impact of different similarity score of different clone detectors in finding co-change candidates in our future studies. Besides this, we want to include some other clone detection tools of different detection mechanism (i.e., tree/ token/ text-based) and subject systems written in some different programming languages (i.e. C/ C++, C#, Python) for extending our research.

Acknowledgment

This research is supported by the Natural Sciences and Engineering Research Council of Canada (NSERC), and by a Canada First Research Excellence Fund (CFREF) grant coordinated by the Global Institute for Food Security (GIFS).

References

- [1] C. K. Roy, J. R. Cordy, R. Koschke, Comparison and evaluation of code clone detection techniques and tools: A qualitative approach, SCIENCE OF COMPUTER PROGRAMMING (2009).

- 535 [2] D. Rattan, R. K. Bhatia, M. Singh, Software clone detection: A systematic review., *Information Software Technology* 55 (2013) 1165–1199.
- [3] C. K. Roy, J. R. Cordy, Benchmarks for software clone detection: A ten-year retrospective, in: *Proc. SANER*, 2018, pp. 26–37. doi:10.1109/SANER.2018.8330194.
- 540 [4] R. S. Arnold, *Software Change Impact Analysis*, IEEE Computer Society Press, Los Alamitos, CA, USA, 1996.
- [5] M. Mondal, B. Roy, C. K. Roy, K. A. Schneider, Associating code clones with association rules for change impact analysis, in: *Proc. SANER*, 2020, p. 11pp.
- 545 [6] T. Rolfsnes, S. Di Alesio, R. Behjati, L. Moonen, D. W. Binkley, Generalizing the analysis of evolutionary coupling for software change impact analysis, in: *2016 IEEE 23rd International Conference on Software Analysis, Evolution, and Reengineering (SANER)*, volume 1, 2016, pp. 201–212.
- [7] J. Svajlenko, C. K. Roy, Cloneworks: A fast and flexible large-scale near-miss clone detection tool, in: *2017 IEEE/ACM 39th International Conference on Software Engineering Companion (ICSE-C)*, 2017, pp. 177–179.
- 550 [8] M. Mondal, C. K. Roy, K. A. Schneider, An empirical study on change recommendation, in: *Proc. CASCON, CASCON '15*, IBM Corp., Riverton, NJ, USA, 2015, pp. 141–150.
- 555 [9] M. Mondal, C. Roy, K. Schneider, Connectivity of co-changed method groups: a case study on open source systems, 2012, pp. 205–219.
- [10] M. Mondal, C. K. Roy, K. A. Schneider, Prediction and ranking of co-change candidates for clones, in: *Proc. MSR 2014*, ACM, New York, NY, USA, 2014, pp. 32–41. doi:10.1145/2597073.2597104.
- 560 [11] M. Mondal, B. Roy, C. K. Roy, K. A. Schneider, Investigating context adaptation bugs in code clones, in: *Proc. ICSME*, 2019, pp. 157–168. doi:10.1109/ICSME.2019.00026.

- [12] J. F. Islam, M. Mondal, C. K. Roy, K. A. Schneider, Comparing bug replication in regular and micro code clones, in: Proc. ICPC, 2019, pp. 81–92. doi:10.1109/ICPC.2019.00022.
- [13] J. F. Islam, M. Mondal, C. K. Roy, A comparative study of software bugs in micro-clones and regular code clones, in: Proc. SANER, 2019, pp. 73–83.
- [14] M. Nadim, M. Mondal, C. K. Roy, Evaluating performance of clone detection tools in detecting cloned cochange candidates, in: 2020 IEEE 14th International Workshop on Software Clones (IWSC), 2020, pp. 15–21.
- [15] J. Svajlenko, C. K. Roy, Evaluating modern clone detection tools, in: Proc. ICSME, 2014, pp. 321–330. doi:10.1109/ICSME.2014.54.
- [16] S. Bellon, R. Koschke, G. Antoniol, J. Krinke, E. Merlo, Comparison and evaluation of clone detection tools, IEEE Transactions on Software Engineering 33 (2007) 577–591. doi:10.1109/TSE.2007.70725.
- [17] C. Roy, J. Cordy, Scenario-based comparison of clone detection techniques, 2008, pp. 153–162. doi:10.1109/ICPC.2008.42.
- [18] E. Burd, J. Bailey, Evaluating clone detection tools for use during preventative maintenance, in: Proc. SCAM, 2002, pp. 36–43. doi:10.1109/SCAM.2002.1134103.
- [19] F. V. Rysselberghe, S. Demeyer, Evaluating clone detection techniques from a refactoring perspective, in: Proc. ASE, 2004, pp. 336–339.
- [20] R. Koschke, R. Falke, P. Frenzel, Clone detection using abstract syntax suffix trees, in: 2006 13th Working Conference on Reverse Engineering, 2006, pp. 253–262. doi:10.1109/WCRE.2006.18.
- [21] S. Ducasse, O. Nierstrasz, M. Rieger, On the effectiveness of clone detection by string matching: Research articles, J. Softw. Maint. Evol. 18 (2006) 37–58. doi:10.1002/smr.v18:1.

- [22] G. M. K. Selim, K. C. Foo, Y. Zou, Enhancing source-based clone detection
590 using intermediate representation, in: 2010 17th Working Conference on
Reverse Engineering, 2010, pp. 227–236.
- [23] T. Software, Tiobe index — tiobe - the software quality company, 2020
(accessed July 6, 2020). URL: <https://www.tiobe.com/tiobe-index/>.
- [24] S. Ducasse, M. Rieger, S. Demeyer, A language independent approach
595 for detecting duplicated code, in: Proceedings IEEE International Confer-
ence on Software Maintenance - 1999 (ICSM'99). 'Software Maintenance
for Business Change' (Cat. No.99CB36360), 1999, pp. 109–118.
- [25] E. Juergens, F. Deissenboeck, B. Hummel, Clonedetective - a workbench
for clone detection research, in: Proc. ICSE, 2009, pp. 603–606. doi:10.
600 1109/ICSE.2009.5070566.
- [26] N. Göde, R. Koschke, Incremental clone detection, in: Proc. CSMR, 2009,
pp. 219–228. doi:10.1109/CSMR.2009.20.
- [27] J. R. Cordy, C. K. Roy, The nicad clone detector, in: Proc. ICPC, 2011,
pp. 219–220. doi:10.1109/ICPC.2011.26.
- [28] M. S. Uddin, C. K. Roy, K. A. Schneider, Simcad: An extensible and faster
605 clone detection tool for large scale software systems, in: Proc. ICPC, 2013,
pp. 236–238. doi:10.1109/ICPC.2013.6613857.
- [29] T. Kamiya, S. Kusumoto, K. Inoue, Ccfinder: a multilinguistic token-based
code clone detection system for large scale source code, IEEE Transactions
610 on Software Engineering 28 (2002) 654–670.
- [30] L. Jiang, G. Mishherghi, Z. Su, S. Glondu, Deckard: Scalable and accurate
tree-based detection of code clones, in: Proc. ICSE, 2007, pp. 96–105.
doi:10.1109/ICSE.2007.30.
- [31] S. Harris, Simian - Similarity Analyser — Duplicate Code Detection for
615 the Enterprise—Overview, 2003 (accessed July 6, 2020). URL: [http://
www.harukizaemon.com/simian/](http://www.harukizaemon.com/simian/).

- [32] C. Ragkhitwetsagul, J. Krinke, D. Clark, Similarity of source code in the presence of pervasive modifications, in: Proc. SCAM, 2016, pp. 117–126. doi:10.1109/SCAM.2016.13.
- 620 [33] T. Wang, M. Harman, Y. Jia, J. Krinke, Searching for better configurations: A rigorous approach to clone evaluation, in: Proc. ESEC/FSE, ACM, New York, NY, USA, 2013, pp. 455–465.
- [34] M. Mondal, M. S. Rahman, R. K. Saha, C. K. Roy, J. Krinke, K. A. Schneider, An empirical study of the impacts of clones in software maintenance, in: Proc. ICPC, 2011, pp. 242–245. doi:10.1109/ICPC.2011.14.
- 625 [35] W. T. Cheung, S. Ryu, S. Kim, Development nature matters: An empirical study of code clones in javascript applications, Empirical Softw. Engg. 21 (2016) 517–564.
- [36] J. Krinke, N. Gold, Y. Jia, D. Binkley, Cloning and copying between gnome projects, in: Proc. MSR, 2010, pp. 98–101. doi:10.1109/MSR.2010.5463290.
- 630