# "EMPLOYABILITY OF 'CODE CLONING' TECHNIQUES IN THE DEVELOPMENT OF LARGE SCALE SOFTWARE SYSTEMS"

**Saksham Rai**

## ABSTRACT

*Large-scale software systems are expensive to build and, are even more expensive to maintain. In PC software, we could possibly have dissimilar sorts of repetition. We ought to note that not every type of redundancy is unsafe. There are diverse types of redundancy in software. Software embodies both programs and information. At particular times redundant is used correspondingly in the sense of unessential in the software engineering works. Sometimes, developers take easier way of implementation by copying some fragments of the existing programs & use that code in their effort. This kind of work is called code cloning. Redundant code is as well recurrently misleadingly permitted as cloned code despite the fact that its point towards that one bit of code which is possibly imitative from the additional one in the original meaning of this particular word. Even though cloning stimuli to redundant code, not each specific redundant code is a specific clone.*

## I. INTRODUCTION

A code clone is a code helping in source documents that is equivalent to or like another [1]. Duplication of code happens every now and again during the improvement of huge programming frameworks. Code cloning is a system of programming reuse and exists in pretty much every product venture. This specially appointed type of reuse comprises in duplicating, and in the end altering, a square of existing code that applies a bit of required usefulness. The primary issue with code clone is coupled uniquely with their comparative code that is in a roundabout way as opposed to straightforwardly which makes it dangerous to distinguish them. The code quality decays, and adjustment turns out to be progressively costly and blunder inclined.

Replicated squares are called clones and the demonstration of copying, tallying slight adjustments, is said cloning. Cloning fundamentally happens in light of the fact that developers find that it is less expensive and speedier to utilize the reorder include then composition the code without any preparation. At times developers purpose on applying new usefulness discover specific working code that finishes a calculation about coordinating to the one wanted to duplicate it all together and then change in place [2] while this is viewed as a major issue in mechanical programming [3]. The nature of code examination, infection perceive, aspect mining, and mistake presentation are different assignments which need the information of linguistically checked code part to encourage code identification significance for programming discovery device investigation.
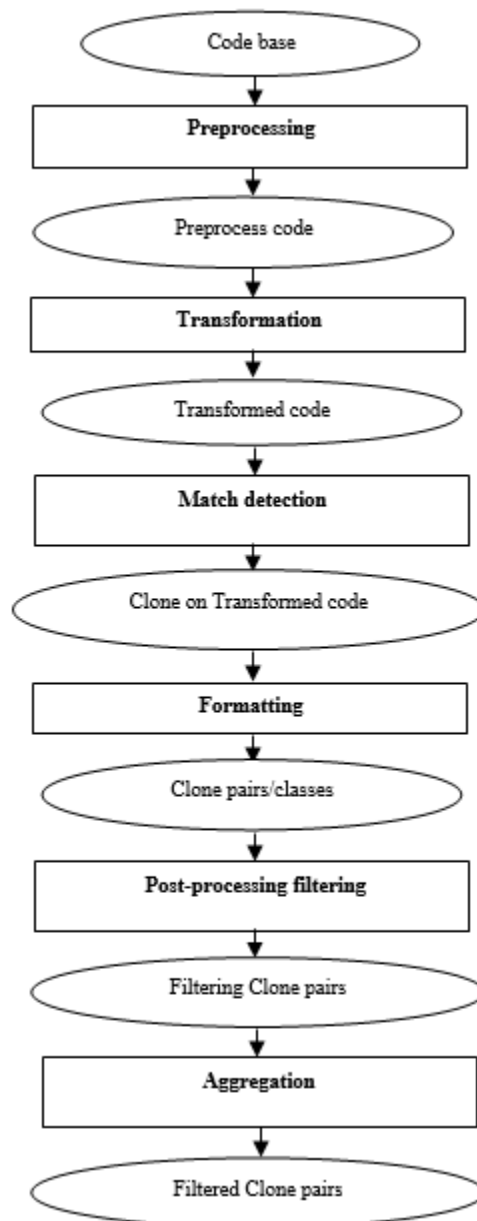
Fig.1: Generic Clone Detection Process

A. Pre-processing

Toward the start of any clone identification approach, the source code is apportioned, and the area of the correlation is resolved. There are three primary targets in this stage:

Remove uninteresting parts: All the source code uninteresting to the examination eliminate is sifted in this stage.

Determine source units: Subsequent to evacuating the uninteresting code, the rest of the source code is apportioned into a lot of disjoint sections called source units. These units are the biggest source parts that might be associated with direct clone relations with one another.

2

Determine comparison units/granularity: Source unit's may be additionally divided into littler units relying upon the correlation procedure utilized by the device. For instance, source units might be subdivided into lines or even tokens for correlation.

B. Transformation

When the units of the examination are resolved, if the correlation strategy is other than printed, the source code of the correlation units is changed to a proper middle portrayal for examination. This change of the source code into a middle of the road portrayal is regularly called extraction in the figuring out the network.

Extraction: Extraction changes source code to the structure appropriate as a contribution to the genuine examination calculation. Contingent upon the instrument, it ordinarily includes at least one of the accompanying advances.

Tokenization: If there should be an occurrence of token-based methodologies, each line of the source is isolated into tokens as indicated by the lexical standards of the programming language of intrigue — the tokens of lines or documents at that point structure the token successions to be thought about.

Parsing: If there should arise an occurrence of syntactic methodologies, the whole source code base is parsed to fabricate a parse tree or (perhaps commented on) theoretical grammar tree (AST).

Control and Data Flow Analysis: Semantics-mindful methodologies produce program reliance diagrams (PDGs) from the source code. The hubs of a PDG speak to the announcements and states of a program, while edges speak to control and information conditions.

C. Normalization

Standardization is a discretionary advance proposed to dispense with shallow contrasts, for example, contrasts in whitespace, remarking, organizing or identifier names.

Expulsion of whitespace: Almost all methodologies ignore whitespace, despite the fact that line-based methodologies hold line breaks. A few measurements based methodologies anyway use arranging and design as a major aspect of their examination.

Expulsion of remarks: Most methodologies evacuate and overlook remarks in the real examination.

Normalizing identifiers: Most methodologies apply identifier standardization before the examination so as to distinguish parametric Type-2 clones.

Pretty-printing of source code: Pretty printing is a basic method for redesigning the source code to a standard structure that expels contrasts in design and separating.

Auxiliary changes: Other changes might be connected that really change the structure of the code, with the goal that minor varieties of the equivalent syntactic structure might be treated as comparative.

D. Match Detection

The changed code is then bolstered into a correlation calculation where changed examination units are contrasted with one another to discover matches. Frequently adjoining comparative examination units are joined to shape bigger units.

E. Formatting

In this stage, the clone pair list for the changed code acquired by the correlation calculation is changed over to a comparing clone pair list for the first codebase. Source directions of each clone pair acquired in the correlation stage are mapped to their situations in the first source records.

F. Post-processing / Filtering

In this stage, clones are positioned or separated utilizing manual investigation or robotized heuristics.

Manual Analysis: After extricating the first source code, clones are exposed to a manual investigation where false-positive clones or fake clones [72] are sifted through by a human master.

Robotized Heuristics: Often, heuristics can be characterized dependent on length, assorted variety, recurrence, or different qualities of clones so as to rank or sift through clone competitors naturally.

G. Total

While a few instruments straightforwardly recognize clone classes, most return just clone matches as the outcome. So as to decrease the measure of information, perform resulting examinations or assemble outline insights, clones might be accumulated into clone classes.

## II. RELATED WORK

ImanKeivanloo et.al (2011) [4] introduced ongoing code clone search was a rising group of clone location looks into that objectives at discovering clone sets coordinating an info code part in divisions of a second. For these systems to meet genuine certifiable necessities, they were to be adaptable and give a short reaction time. Our examination exhibited a half breed clone search approach utilizing source code example ordering, data recovery bunching, and Semantic Web thinking to separately accomplish less reaction time, handle false positives, and bolster computerized gathering/questioning. Yang Yuan et.al (2011) [5] introduced CMCD; a Count Matrix based method to detect clones in program code. The key model behind CMCD was Count Matrix, which was created while counting the occurrence frequencies of every

variable in conditions specified by pre-determined counting situation. Because the characteristics of the count matrix do not change due to variable name replacements or even switching of statements, CMCD works well on various hard-to-detect code clones, such as exchange statements or deleting a few lines, which are difficult for other state-of-the-art detection techniques. DebarshiChatterji et.al (2011) [6] described a study that investigates the usefulness of code clone information for performing a bug localization task. In this study 43 graduate students were practical while identifying defects in both cloned & non-cloned portions of code. The goal of the study was to understand how those developers used clone information to perform this task.

## III. CLONE TYPES

Code clone could be of any sort that all rely on upon the developer's method & ability of utilizing the code which varies from replicating as it is to matching the code however with some change which would be done at diverse level in the method. In software scheme code pieces predominantly demonstrates two sorts of similarities. One clone type of similarity considers textual similarity), & another second considers the semantic level that the clone code essential to have the identical behaviors, means the functional similarity.

A. Textual Similarity: Two code fragments can be similarly based on the similarity of their program text we differentiate the subsequent sorts of clones. The following types of clones are discussed in order to find textual similarity [2].

1) Type I: In Type I clone, a copied code fragment is the same as the original. However, there might be particular variations in whitespace (blanks, new line(s), tabs etc.), remarks and/or designs. Type I is widely known as exact clones

2) Type II: A Type II clone is a code fragment that is the same as the original except for some possible differences about the corresponding designations of user-defined identifiers (name of variables. constants, class. methods & so on) layout, identifiers, remarks, literals, & sorts. The specific reserved words & the sentence structures are essentially the same as the original one.

3) Type III: Type DI is copied with further modifications. E.g. a new statement can be added, or some statements can be detached along with various dissimilarities in layout, identifiers, remarks, literals, & sorts. The structure of code fragment may be changed & they may even look or behave slightly differently. This kind of clone is hard to be discovered, for the reason that the whole framework understanding is needed.

B. Functional Similarity: Two code fragments can be similarly based on the similarity of their functionalities without being textually similar. If the functionalities of the two code fragments are identical or similar i.e., they have comparable pre as well as post circumstances referred as Type IV clones or semantic clones [8].

4) Type IV: Type IV clones are the consequences of semantic similarity between two or extra code fragments which could accomplish the same computation however actualized through

diverse syntactic variations. In this category of specific clones, the cloned part is not necessarily copied from the first one. Two code fragments may possibly be established by two different programmers too.

## IV. DETECTION OF CLONE AND USES

Clone detection is useful in finding malicious software [8]. By comparing, one malicious software with another, it is possible to find the matched parts of one software system with another. Some applications of clone detection are as follows:

a) Plagiarism Detection: In Projects Plagiarism Detection is one of the firmly related regions of clone recognition [11]. Clone identification procedures can be utilized in the area of literary theft location. A clone location instrument, for example, token-based CC Finder has been connected in distinguishing unoriginality.

b) Copyright Infringement: In Projects Plagiarism Detection is one of the firmly related regions of clone recognition [11]. Clone identification procedures can be utilized in the area of literary theft location. A clone location instrument, for example, token-based CC Finder has been connected in distinguishing unoriginality. [9].

c) Clone Detection: in Models, Clone detection is also used in models [10]. The phenomenon is not restricted to code but usually occurs in models in a very similar way. So it is likely that model clones are as unfavorable to model quality as they are to code quality. Clone detection is also used in the data flow model. General Model. LIE Domain Model [10, 11].
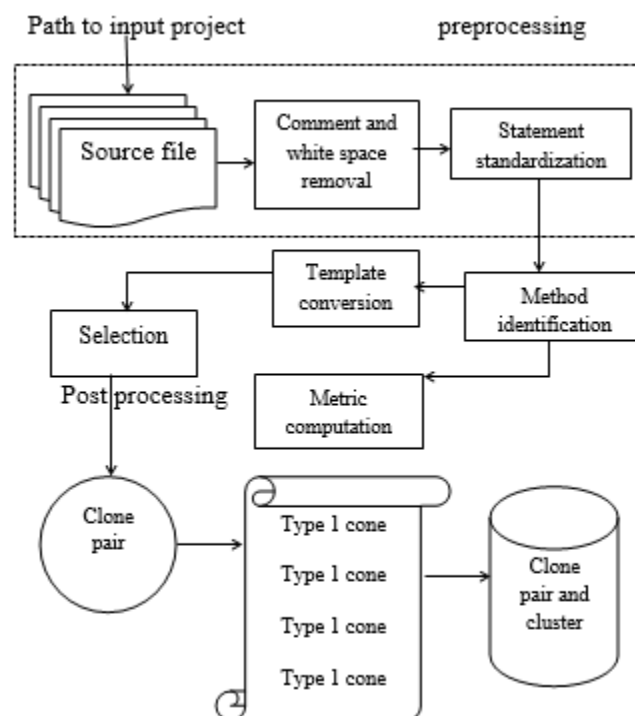


Fig.2: Schematic diagram of Clone Manager [14]

# V. DETECTION CLONE TECHNIQUES

Clone Detection is a dynamic research territory since 1990 "s [6]. Code clone location is exactly identified by the upkeep of software, code overhauling & along these lines making the code more effective. The impression on clone identification demonstrates the dissimilar strategies & calculations to distinguish clones [3]. Clone position methodologies are comprehensively considered into five procedures which are portrayed underneath:

A. Text-Baseded Technique

One of the quickest clone recognition draws near. It can undoubtedly manage type 1 clone and with extra information change, the sort two can likewise be taken consideration. The more up to date message-based clone location procedure that depends on spot plots. A dab plot is a two-dimensional diagram where the two tomahawks rundown source elements. In this methodology, the lines of a program are examination elements. In the event that x and y are equivalent, there is a dab at arrange (x, y). Two lines are viewed as equivalent on the off chance that they have similar hash esteem. Dab plots can be utilized to envision clone data; diagonals in speck plots are recognized as a clone. The location of clones in dab plots can be computerized, and string-based powerful example, coordinating is utilized to look at entire lines. Diagonals that have holes demonstrate type 3 clones.

B. Token-Based Technique

Token-based system is like content-based strategy. Be that as it may, rather than accepting a line of code as portrayal straightforwardly, lexical analyser changes over each line of code into a grouping of Token [12]. After information esteems and identifier are substituted by some uncommon tokens. The token groupings of lines are thought about effectively through a postfix tree calculation. The outcome is likewise exhibited in speck plot chart. This system is to some degree slower than the content-based method on account of the tokenization step. Notwithstanding, applying postfix tree coordinating calculation the time multifaceted nature is comparative as a content-based strategy. By breaking the line into tokens, it can without much of a stretch recognize both sort 1 and type 2 clone, with token channel connected. The aftereffect of the clone can be controlled absolutely, for example, avoid any uninterested data.

C. Abstract Syntax Tree (AST) - Based Technique

In AST based method the program is parsed into a parse tree or dynamic punctuation (AST) with a parser of the language of intrigue. At that point, utilizing a tree coordinating technique, comparable subtrees are inspected in the tree. At the point when a match is discovered relating, so code of the same subtrees is returned as clone couples or clone classes. The data is accessible in the parse tree of AST. The variable names and strict worth the source code is disposed of during the tree portrayal: still, it is conceivable to utilize increasingly advanced clone

7

recognition instruments. By utilizing AST as code portrayal gives this method a superior comprehension of the framework structure. In any case, the parsing source document is as yet an over the expensive top procedure on both time and memory.

# VI. CONCLUSION

Code clone is a big problem. A copy and paste activity which is done by the programmer is the main reason of code cloning. It looks like a simple and effective method, these copy and paste activities are not documented. Which create a bad effect on the software quality and duplication also increase the bug probability and maintenance problem? Cloning of code has become one of the easiest ways to complete a project, who does not want to invest their time on doing programming their project. It's a loss for those who really work hard for the project coding. The date no such method has present who can evaluate the cloning for several languages with one piece of code. In this paper also present the different types of code clone techniques like text-based Token-based comparison and Abstract Syntax Tree-Based Comparison. In this paper describe how to effective these techniques in code clone.