# EMPLOYABILITY OF ABSTRACT SYNTAX TREES IN THE EFFICACIOUS REPLICATION OF PAID OR PREMIUM SOFTWARE AND ITS EXTENSION TO MULTIPLE FUNCTIONALITIES – STRING MATCHING AND FUNCTIONALITY BASIS

**Lavaneesh Sharma**

## ABSTRACT

*Code clone is the technique in which developer produces an exact replica of an original product or software. Extending upon the previous statement, it is a technique which requires no sort of parent to as to actually create a software. A software clone, typically doesn't require any original source code for functioning of the software, and herein this thesis we will try to as accurately replicate the software that is already extant in paid or premium version format. The AST is utilized to detect clone and there exist many method and using them we can detect the software clone like string matching, functionality basis.*

## INTRODUCTION

Code replication or recreating a code segment and imminently reuse it by staying using it or without using it any progressions is a striking code peculiarity in programming upkeep. A couple of examinations demonstrate that five to twenty percent of a product's structure can contain a replicated program or code, which is essentially the outcome of copying parent code in parts and using them by coalescing it with our required program in part or none, at all. One, of the genuine insufficiencies of such replicated segments arises if there has, been a detection of any sort of programming bug; the different pieces like it (the one under consideration) should be examined thoroughly to detect the possible nearness of a comparable bug in the tantamount areas. Re-development of the copied code is yet one of the upmost problems in programming support but a couple of examinations ensure that re-incorporation of explicit clones isn't fascinating and involves a risk of ousting them. Any case be, it is furthermore majorly agreed that clones should in any occasion be recognized.

Code cloning is notable issue in programming and therefore leads to poor programming quality tasks. The motivations to duplicate code sections are:

430

1. Making a duplicate of code is straightforward and quick, instead of composing from scratch
2. Producing more source code prompts better impetuses for software engineers in industry.
Tools and techniques for distinguishing copied code are of primal importance in programming upkeep look into. A segment of the definitions identified with code-cloning are talked about underneath:

*Code Fragment*: A code section (C.F.) is any arrangement of lines of code (with or without remarks). It very well may be of any size, e.g., that defines a function, succession of proclamations or a starting-ending chunk.

*Code Clone*: A code clone is an important part of source records, comparable or indistinguishable from another code segment. A code divide (*CPi*) is a copy of other code segment (*CPj*), on the fact that they have some connection, *iff*,
$f(CPi) = f(CPj)$
where *f* is any likelihood function.

*Clone Pair*: A complementary clone segment is said to be uniquely defined if there can be an identifiable clone relation amongst them, i.e. a clone pair.

Code clones (or copies) are a commonplace in all contemporary products. Customarily, code clones have been viewed as an issue and are generally unsuitable towards quality aspect of programming. Eventually, it was previously understood that despite the fact that code clones are a potential support overhead and source of bugs, they can't generally be avoided and their expulsion may not generally be conceivable. Therefore, best case scenario they are viewed as code scents and they should be overseen inside source, code in an increasingly viable manner.

The present arrangement of tools that endeavour to redress the zone of code clone management extend from giving recognition to and subsequently, following development of clones. There are a few devices likewise, giving recommendations to refactoring. Anyways, the consequences of these tools and corresponding solutions towards fixing problematic clones are manageable when the original code materialises to be small (i.e. code lines). As code's base builds, the aftereffects of a clone; a detection tool can be colossal; making it exceedingly hard to extricate the critical clones that should be expelled or fixed. It is historically confirmed by numerous analysts that almost five to twenty percent of large software frameworks are clones, hence approving the purposes behind countless occasions of code clones in the reports. Isolating the code clones of significance from mammoth number of clones is a noteworthy issue as a result of the complementary reasons:

431

-Real-world programming ventures quite often run shy of time & assets. There is a need to accomplish within the projects constraints such as time and feasibility, and furthermore, simultaneously fix the issues recognized by code assessment and clone detection devices before each release.
-Given these requirements, if there should arise an occurrence of code bases of enormous sizes, fixing all clones in such a situation is certainly not a reasonable arrangement.
-Each task has explicit needs and necessities and no venture will ever need to fix every single instance of code clones, detailed by devices that identify clones.

None of the extant tools have as of yet, endeavoured to necessitate the accommodation of these clones as long as the request goes by, where the remedial action might, be taken. Furthermore, scant of those have endeavoured to provide a unifying structure towards fixation of code clones. What is often usually required, for clone organization, is the device that can organize the clones for fixing and if conceivable, give insights towards fixing.
Likewise, such a strategy or criteria, must be useful to all modern programming, regardless of size, programming language, or area. Heading this way, there is absenteeism of all encapsulating perspective or system for addressing code clones and their compelling administration.
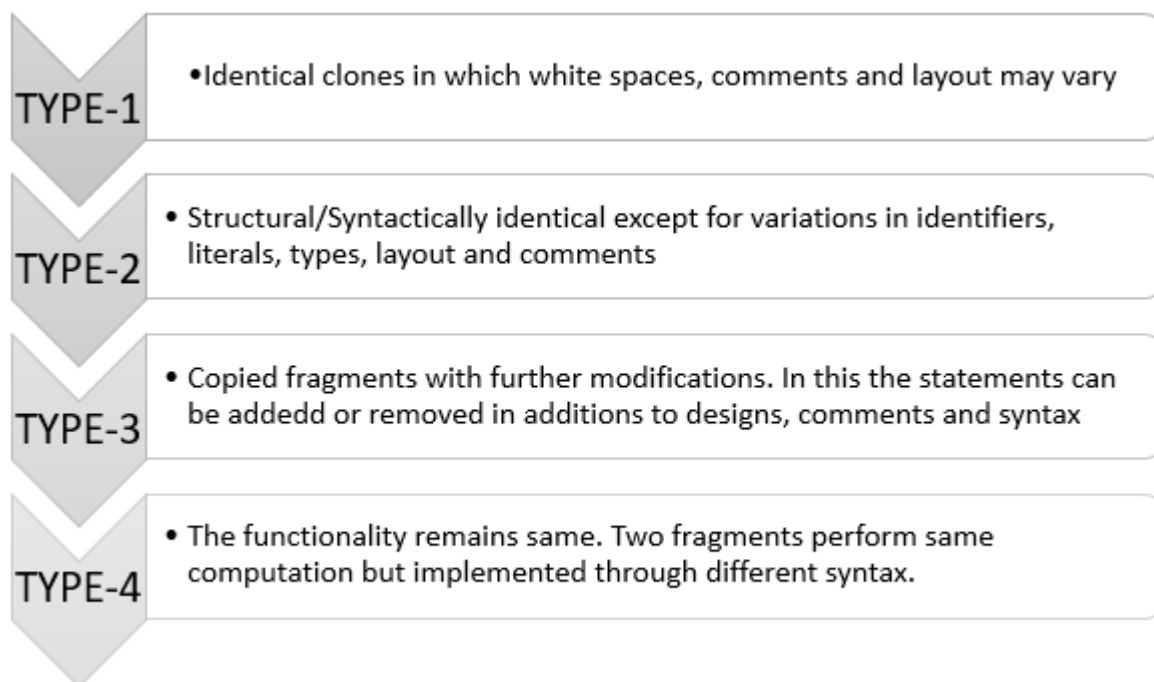
The proposed techniques here, attempts by helping the undertaking of clone management by helping a software engineer to house the needs to sets of clones recorded in the code clone reports. The developer can join different elements (clones effect additionally to support, quality perspectives, and refactoring goes) to settle on which clones are progressively more critical and necessarily be prioritized and those that are of lower need. It supports to choose, what we term as, the Programming Clone Quality – quality | of code segments recognized as clones, at par with its adherence to programming guidelines. In view of the measurements and qualities related with these variables, the conceivable fixes for clones can be derived.

## TYPES OF CLONES:

Code clone recognition process is defined to be the reprocessing of the original code. These clones make the entire code redundant. The code cloning additionally prompts the additional bugs in program. The code cloning majorly affects the product business as it edifices the plan of the product venture and furthermore it is tough to mend the framework. Source code cloning asserts itself as a grave threat to the practicality of a product framework. The methods chiefly utilized for
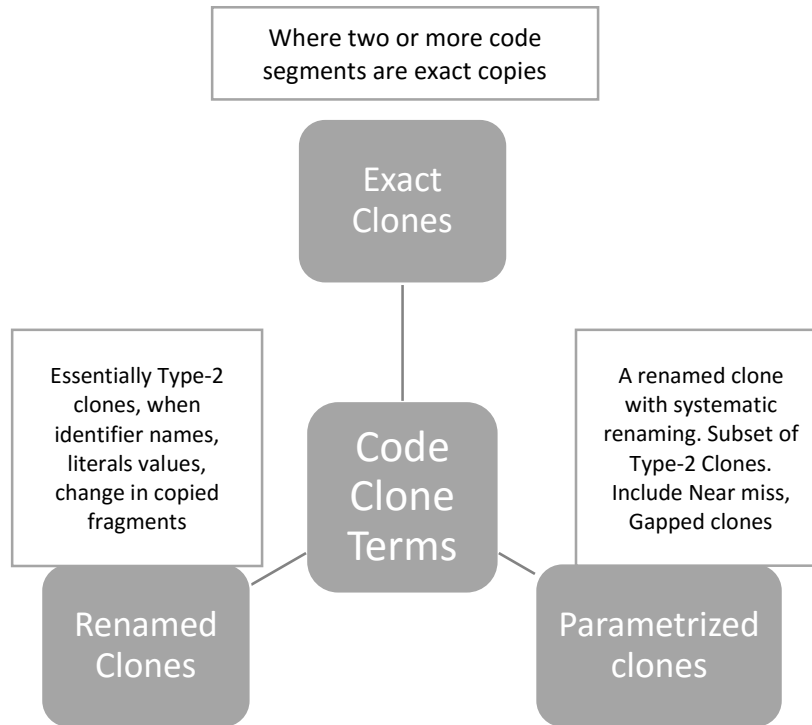
432

code clone location are: Text-based, Token-based, Tree-based, Program reliance chart based, Metric-based. In paper we survey the diverse system(s) for recognizing the code clones.

Code cloning means adapting a few portions, variables, functions and afterwards pasting them in another program is code cloning. In each product nearly seven to twenty-three percent code is duplicated. This entails, the product redundant and it's support cost increases. Particularly many open source code are generally duplicated. So discovery of these clones is exceptionally vital to keep up the product cost. The different discovery strategies are utilized on premise of which sort of code clone are inherently associated with the product.

**TYPE-1**
- Identical clones in which white spaces, comments and layout may vary

**TYPE-2**
- Structural/Syntactically identical except for variations in identifiers, literals, types, layout and comments

**TYPE-3**
- Copied fragments with further modifications. In this the statements can be addedd or removed in additions to designs, comments and syntax

**TYPE-4**
- The functionality remains same. Two fragments perform same computation but implemented through different syntax.

The Type-1, Type-2, Type-3 code clones are known as syntax clones whereas Type-4 is called semantic clone.

- Clone Pair: If there is an equivalence relation between two code segments, then they form a clone pair.
- Clone Class: It is defined as collection of similar code segments. Each code segment in a clone class form a clone pair with other code segment of that class.

433

Where two or more code segments are exact copies

Exact Clones

Essentially Type-2 clones, when identifier names, literals values, change in copied fragments

Code Clone Terms

A renamed clone with systematic renaming. Subset of Type-2 Clones. Include Near miss, Gapped clones

Renamed Clones

Parametrized clones

## PARAMETERS USED FOR CODE CLONE DETECTION METHODS

There are a few code clone detection methods. The correlation of these structures is quite much, which necessitates to prize the exact process for the issue proclamation. Various parameters have been decided for given comparison via five procedures (content based, token-based, tree-based, diagram based, metric-based). These parameters are otherwise called the clone challenges. A portion of the parameters are expressed underneath: -

| | |
|---|---|
| Portabiltiy | • The approach(s) must be convenient for different dialects and lingos as many programming languages are incorporated with several dialects. It is standard, that clone identification procedures, be effectively compact and configurable for various dialects |
| Precision | • The methods ought to be sound enough to make a division among the less number of false positives. It is likewise being said that strategies should discover the copied code with enhanced precision. |
| Scalability | • It is tough to unearth the clones of code from the huge and complex framework. The strategies would be versatile to effectively manage the huge and complex framework with efficacious utilization of memory |
| Robustness | • A nice framework should be lively for different changing activities that are associated on copied piece and recognize the clones with higher precision and audit. |

434

# PROBLEM STATEMENT

Our goal is to develop a program to find out the exact location of code clones and their types. There are numerous approaches involved in checking code clones, but we are going to use Abstract Syntax Tree (AST). We implemented String Distance currently to label the clones as Type1, Type2 or Type 3. Categorization is done via the master expression strings.

At the point when programming developers actualize finishing the program, comparable code section may show up in various, pieces of the product framework. Seldom, it is done purposefully and in some cases clone codes happen inadvertently. More often than not, codes are reused by the developers for their outstanding arrangements. At the point when clone code is deliberately presented in programming program, the developer may accomplish time productivity, yet cedes control of the executed programming framework. Aside from previously mentioned issue there is few more issues that can happen because of code cloning. Issues are:

-Bug transmission probability is enhanced, as one segment is copied from other.
-New Bugs can be introduced if the developer is not prudent in his programming.
-Lack of good inheritance framework and abstraction may lead to poor design and robustness.
-Maintenance is enhanced. Even presence of one bug calls for the thorough enquiry of the entire code. Duplication via maintenance is deleterious.

# PROPOSED APPROACH AND IMPLEMENTATION

Contrast each subtree with other subtree in AST. This can be done in $O(N^3)$ where N is the number of hubs in AST.

Implement a hash function to compare the different subtrees

If the likelihood of hash function is similar or passes a threshold (weight of subtree passes the edge weights), then they are identifies as clones

## BASIC ALGORITHM

```
Clones = NULL
For each subtree i:
  If (weight(i) >= Threshold)
   Then (hash i to bucket),
     For each subtree i and j in the same bucket:
       If (Compare tree (i,h) > Similar Threshold),
       Then,
        For each subtrees of i:
          If (IsMember r(clone s, s),
           Then Remove Clone Pair (Clone s,s),
       For each subtree s of j:
        If (IsMember r (clones s,s))
        Then Remove Clone Pair (Clone s,s)
  Add Clone Pair(Clones i,j);
```

## SEQUENCE DETECTION ALGORITHM

The list 's' is created describing sequences

For k = MinimumSequenceLengthThreshold to MaximumLengthThreshold

Place all subsequences of length k into buckets according to subsequence hash

For each subsequnce i and j in the same bucket

If Compare Sequences(i,j,k) > Similarity Threshold, Then

{RemoveSequenceSubclonesOf (clones i,j,k) AddSequenceClonePair (clones i,j,k)}

436

## IMPLEMENTATION

1.  We have instigated a python script, where we enter the title of python script to be comment cleaned first, as comments repeat themselves quite often and they would appear to a great extent in the clones, goal is to concentrate on code rather than comments. File to be comment cleaned: test1.py

We use tokenization and lexical analysis to remove the comments.



2.  After this a file is generated "cleanedlib.py" in which the comments were removed.



**Fig 1: Before Comment cleaning**

437

**Fig 2: After Comment cleaning**

3.  Then the cleaned file is passed to the script that generates the AST and look for clones, we store these clones as a dictionary and get the list of clones.



4.  As there are multiple lines in the code, where clone is detected, we count those and put those back in a separate file and csv to easily analyze. Types of clones are also stored separately.

438

5. Clone digger a famous and accessible device, is used to contrast this process It can't separate the clones, despite being one of the major devices; uses string division for separation. Finally, eclectic strategies exist to characterize the utilizing just direct string coordinating or measurements, which happen to be not dynamic and don't give exact examination



The clone prioritization framework takes a shot at a chose rendition of code base acquired from the adaptation control framework, on which clone identification apparatuses are run, coming about

439

in (now and again) huge reports comprising data in accordance with the clones in that peculiar form. This framework which is planned for supporting clone management utilizes results obtained because of running different investigation devices like – Static examination instruments, Refactoring workbench, Program perception apparatuses, Metric devices – aftereffects of whom are mapped to a chose quality model. We have acquired the EMISQ quality model (which depends on ISO 9126 standard) for our execution. As it were, it decides the Software Clone Quality – the amount to which practicality is influenced and the consequence on the measured code quality depends to amount of infringement and its seriousness. It additionally contemplates the impacts of refactoring the clone.

So as to settle on which, set of clones should be later surveyed up on need, a few benchmarks should be deliberated. Drawing from our past encounters in code quality and refactoring surveys, we have devised tactics to consequently distinguish basic clones on premise of different elements. These variables help decide the request for the clone results. Each factor comprises of a few estimates which in total decide the weight or the need an incentive for the clone class and for each clone individually in the clone class.

## CLONE DETECTION USING ASTs



**Architecture Diagram**

As an initial phase in the clone discovery process is the parsing of the original code and an AST is delivered. Three fundamental calculations are required to discover clones. The premise of the primary calculation is to identify sub-tree clones. The subsequent one is grouping recognition calculation, which is concerned about the detection of variable-size arrangements of sub-tree clones. It is additionally utilized to distinguish articulation and assertion succession clones. The third calculation center in ever more intricate, close miss clones by looking to sum up different clones. The subsequent distinguished clones would later be able to be duplicated.

# CLONE REMOVAL



Code clone removing is organized in two parts. In the primary stage, a point by point investigation of recognition of code clones is done by also including the dynamic language structure tree. This clone identification presents straightforward and useful techniques for finding careful and close miss clones over self-assertive program sections in source code by analyzing unique language structure trees. In the subsequent stage, we center around how the consequences of stage one can be introduced so as to direct an intelligent refactoring/clone evacuation process.

# ABSTRACT SYNTAX TREES

A AST address the major components of a grammar of a programming language, like punctuation trees that etymologists use for human dialects. The tree centers around the guidelines instead of components like props or semicolons that end proclamations in certain dialects. The tree is various leveled, with the programming components articulations, separated into their parts. For instance, a tree for a restrictive articulation has rules for factors hanging down from the required administrator.

ASTs are generally utilized in compilers to check code for exactness. On the off chance that the created tree contains mistakes, the compiler prints out an error message. ASTs are utilized in light of the fact that a few software's can't be spoken to in a setting free sentence structure, for example, certain composition. ASTs are profoundly explicit to programming dialects, yet research is in progress on all-inclusive sentence structure trees.

In programming building, a hypothetical AST, or just accentuation tree, is a true representation of the extraordinary syntactic structure, of source code programmed in specific language. Each center point of the tree demonstrates an advancement that can or cannot occur in the source code. The language structure is "hypothetical" as it does not completely specify the punctuation. For instance, closed grouped sections are encapsulated in tree structure and a syntactic form like a happening, that condition clarification may be implied by techniques for a single center point with three branches. This perceives hypothetical phonetic structure trees from solid language structure trees, by and large doled out parse trees, which are ordinarily worked by a parser during the source code translation and totaling process. Whenever made, supplementary information is added to AST by techniques for results preparation, e.g., consistent assessment. AST are data structures commonly used in compilers to address the structure of program code. An AST is commonly the consequence of the language investigation time of a compiler. It normally fills in as a moderate depiction of the program through, a couple of stages that the compiler needs, and unequivocally influences the last output of the, compiler.

AST is helpful as:
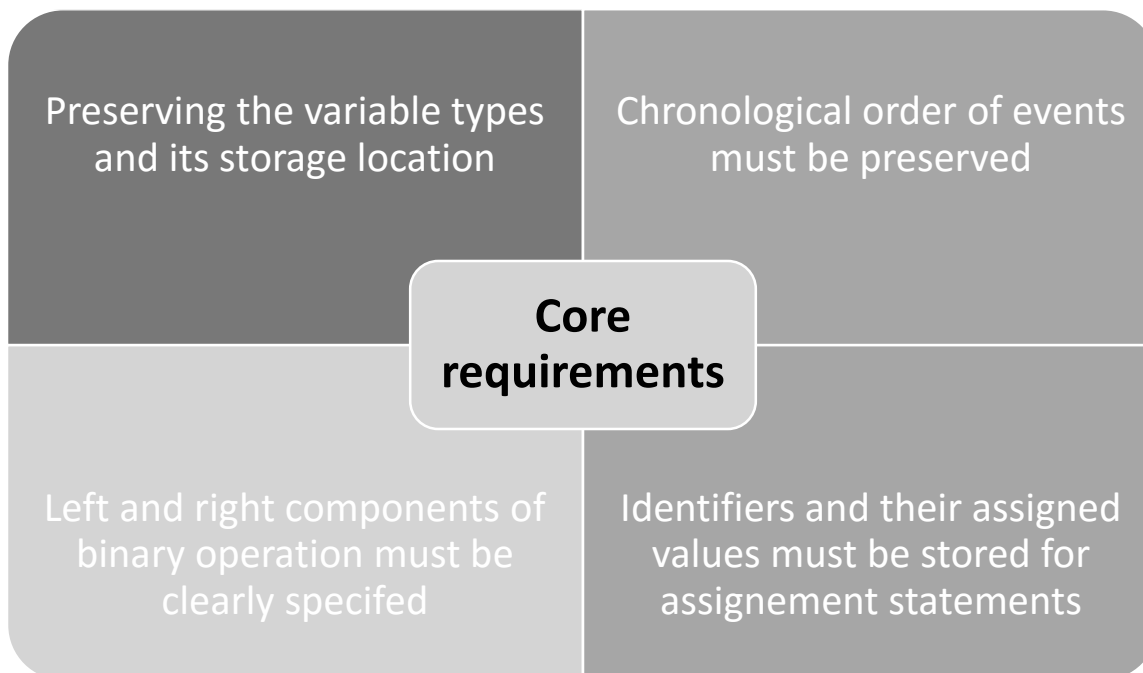-AST devalues redundancies like delimiters and punctuations.
-Useful data is added as a feedback into an AST, due to successive analysis by the compiler.
-AST alters and gives a dimension to contain new information of every encapsulated element, otherwise hard to obtain with a program's source code.

ASTs are required in light of the natural thought of programming lingos and their documentation.

Dialects are mostly flawed. In order to dodge this dubiousness, programming dialects are consistently demonstrated as a Context-Free Grammar(CFG). Regardless, there are normally parts of programming vernaculars that a CFG can't express, anyway are a bits of language and are recorded in its assurance. These are nuances that require a setting to choose their authenticity and lead. As an example, if a language empowers new sorts to be declared, a CFG can't envision the names of such sorts nor the way by in which they are allowed to be used. Despite whether a language has an already well established set of laws, maintaining a real use generally requires some unique kind of circumstance. Another model is duck composition, where the sort of a segment can change dependent upon setting. Admin over-troubling is one more instance, where right usage and last limit are settled subject to the particular instance. Java gives an eminent model, wherein the symbol ('+') is both numerical extension and association of strings.

Actually, there are other data structures connected with the interior activities of a compiler, the AST plays out a unique role. During the initial arrangement, the language structure assessment delivers, a compiler conveys a parse tree. This parse tree can be used to perform for an all- purpose, components of a compiler by strategies intended for syntax composed elucidation. Depite, of the way that this strategy can incite an undeniably successful compiler, it clashes with the item fabricating gauges of forming and keeping up projects.

## DESIGN:

| | |
|---|---|
| Preserving the variable types and its storage location | Chronological order of events must be preserved |
| **Core requirements** | |
| Left and right components of binary operation must be clearly specifed | Identifiers and their assigned values must be stored for assignement statements |

A few tasks will consistently require two components like two terms for expansion. As it so happens, some language builds require a subjectively enormous number of kids, for example, contention records go to program's shell direction. Accordingly, an AST implemented so as to execute a code programmed in given language, likewise needs to be adaptable enough to take to justify for snappy expansion of an obscure amount of offspring's. Another structural necessity for AST is to convert an unparsed AST into source code structure. The source code thus obtained must be like the parent in appearance and unidentifiable whilst running after recompilation.

## CONCLUSION

In this thesis, the issue area is addressed and discussion about the peculiar kinds of code clones is delivered. According to expected focal points, an inferred essentialness of code clone identification is presented and a genuine plan is browsed through various methods, wherein, we have proposed applied accentuation tree based clone ID technique, using hashing capacity, with the ultimate objective that this capacity can hash a declaration to a key which recommended parse tree to absolutely be clear. The hashing limit can hash a similar kind of accounts into a comparable key, irrespective of the wide array of operands they have. Along these lines, we can recognize the clones that are bigger than various frameworks recognized. Lastly of this examination, all targets are cultivated and the criteria met along these lines we can propose that the issue is illuminated.

A pragmatic business technique for recognizing close miss and progression clones, is shown at a scalable level. The approach relies upon assortments of techniques, tools and processes for compiler essential sub expression transfer using hashing. The procedure is clear to realize, via the standard parsing advancement, distinguishes clones in self-emphatic language creates, and process macros that allow the exodus of the clones, without actually impacting the activity of the program. We associated the environment to an authentic usage of moderate scale, and certified past appraisals of clone thickness of 7-15%, proposing that there's, a "manual" programming designing methodology "excess" steady. Robotized systems can recognize and oust such clones, cutting down the estimation of this steady, at accompanying investment funds in program building or upkeep costs. Clone recognition apparatuses additionally have great potential for helping area investigation. As time passes, by we call out for more advance study in this field, basis the grounds of framework we have tried to establish.