*Available Online through*        *Research Article*

**www.ijptonline.com**

# IDENTIFICATION OF STRUCTURAL CLONES VICTIMISATION ASSOCIATION RULE AND CLUMP

**E.Jacob Evanson Solomon* and A.Rama**

Department of Mechanical Engineering, Bharath University, Chennai.

Assistant Professor, Department of Information Technology, Bharath University, Chennai.

Email:rama_j1@yahoo.com

**Abstract**

Code clones square measure similar program structures of considerable size and vital similarity. Easy clone set fashioned by similar code fragments in code. The matter is that the vast variety of straightforward clones generally rumored by clone detection tools. We tend to discovered that revenant patterns of straightforward clones – supposed structural clones - usually indicate the presence of attention-grabbing design-level similarities. we tend to propose a way to sight some specific sorts of structural clones from the recurrent mixtures of co-located easy clone. we discover the patterns of co-occurring clones in numerous files victimisation the frequent item set mining (FIM) technique. Finally, we tend to perform file clump to sight those clusters of extremely similar files that square measure seemingly to contribute to a design-level similarity pattern. we tend to implement the structural clone detection technique in a very tool referred to as CCFinder. Detection of clones provides many advantages in terms of maintenance, program understanding, reengineering and use.

**Keywords:** Style ideas, maintainability, reengineering and reusable code.

## 1 Introduction

CODE clones square measure similar program structures of considerable  size and vital similarity. many studies counsel that the maximum amount as 20-50 p.c of huge code systems incorporates cloned code. Knowing the placement of clones helps in program understanding and maintenance. The detection and future resolution of clones by refactoring, operate calls, macros and templates etc., however, guarantees decrease in maintenance prices and code size. Within the past decade, clone detection and backbone has considerable  attention from the code engineering analysis community and lots of clone detection tools clone detection has been centered on detection similar code fragments – supposed easy clones. we tend to discovered that revenant patterns of straightforward clones usually

indicate the presence of attention-grabbing higher-level similarities that we tend to decision Structural clones whose unification not solely brings a lot of size reduction, however conjointly helps in understanding the planning of the system for higher maintenance and future sweetening. The limitation of considering solely easy clones is thought within the field. the most drawback is that the vast variety of straightforward clones generally rumored by clone detection tools. There are variety of tries to maneuver on the far side the data of straightforward clones. we tend to discovered that at the core of the structural clones, usually there square measure easy clones that exist and relate to every alternative in bound ways in which. We tend to planned a way to sight some specific sorts of structural clones from the recurrent mixtures of colocated easy clones. we tend to enforced the structural clone detection technique in a very tool referred to as CCFinder, enforced in C++. it's its own token-based easy clone detector. Our structural clone detection technique works with the knowledge of straightforward clones, which can come back from any clone detection tool. It solely needs the information of straightforward clone sets and therefore the location of their instances in programs. As structural clones usually represent some domain or style ideas, their information helps in program understanding, and their detection opens new choices for style recovery that square measure each sensible and ascendable.

Representing these recurrent program structures of huge coarseness in a very generic kind conjointly offers attention-grabbing opportunities for use and their detection becomes helpful within the reengineering of inheritance systems for higher maintenance. We are able to notice clone patterns in numerous units of code, either ways or categories or elements or modules, gaining helpful insights into the biological research state of affairs at totally different levels of abstraction. we've got ab initio tried this approach at file level, by finding the oftentimes occurring clone patterns in numerous files and analyzing those patterns, with promising results. detection the oftentimes co-occurring clone categories in numerous files, we are able to isolate the teams of files that have sturdy similarity with one another. this is often achieved by a clump rule that we've got devised for this explicit drawback. These clusters of extremely similar files kind basic structural clones. The rest of this paper is organized as follows: In Section two, we tend to outline sorts of structural clones-higher level similarities in programs. Section three describes our detection structural clones with data processing. Section four describes the implementation of CCfinder. Section five describes a mechanism to make generic illustration of structural clones found within the system for higher maintenance and use. Section half-dozen presents the connected add higher-level similarities and style recovery. Section eight concludes the paper and presents future work.

## 2  Structural Clones- Higher Level Similarities In Programs

We tend to describe intimately the development of upper level similarities, that we tend to decision structural clones. We tend to outline structural clones as similar program structures which will be analyzed hierarchically, at several levels of abstraction, with similar code fragments at all-time low of such hierarchy. Locating these higher level similarities will have vital worth for program understanding, evolution, reuse, and reengineering.

### 2.1  From Easy Clones To Structural Clones

We tend to primarily specialize in similarity patterns representing style ideas or solutions which will be of great importance within the context of understanding, maintaining, reengineering or reusing programs. we tend to use the term structural clone to mean similar program structures that square measure configurations of lower-level similar program entities. Therefore, our structural clones might kind a hierarchy of clones, with cloned code fragments at all-time low level.

### 2.1.1 File-Level Structural Clone

Functions shown within the same shade square measure clones of every alternative (e.g., staff_fn1, task_fn1, project_fn1). The Relationship between the functions is 'same file', that holds between fragments of a similar file, despite the order during which they seem. The 3 host files editStaff.php, editTask.php and editProject.php perform similar tasks, however belong to a few totally different modules (i.e., employees module, Task module, and Project module). Provided these structural clones cowl a considerable portion of the host files, {we can|we will|we square measure able to} think about the 3 files as abstract entities that are clones of every alternative, as mentioned within the previous section. This illustrates however the idea of structural clones helps North American nation to maneuver from smaller entities (in this case, functions) to larger entities (in this case, files). These files will currently be thought-about as entities in forming a better level structure.

### 2.1.2 A Module Level Structural Clone

Clone per the definition, structural clones of upper coarseness is created of structural clones of lower coarseness. for instance, a module-level structural clone will incorporates file-level structural clones. Such a state of affairs is illustrated by the structural clone found in a very internet portal implementation, during this portal, files happiness to every module square measure hold on in a very separate folder. Every module contains a collection of files providing module-specific implementation of bound common functionalities (e.g., create, display, edit, delete). once the module functionalities square measure similar, every of those common files winds up being file-level structural clones of

their counterparts in alternative modules. One such case was the premise for previous example. At a bigger coarseness, the modules employees, Task and Project is thought-about structural clones wherever every structure has four files as entities and therefore the relationship 'same folder' among Note that the module Project doesn't carry a deleteProject.php file. Still, there was enough similarity among Project and alternative modules to think about all of them as structural clones of every alternative.

## 2.1.3 Multiple Structural Clones within The Same File

Multiple Structural Clones within the same file Multiple Structural clones may also exist within the Indian file two shown the four structural clones square measure structures of code fragments that square measure a part of the templates representing various hashed associative containers. every structural clone covers a big a part of the example it belongs to; thus we are able to think about these templates as 'abstract entities' (ignoring their internal structure) and kind a clone category of 4 clones at successive higher level. What is more, if the 2 templates gift in one file square measure joined to every alternative with the 'same file' relationship, we've got a structural clone category of 2 structures, one in every file. Raising the amount of abstraction by an additional step, we tend to discovered that the 2 templates gift in every file cowl the files considerably, therefore the files may also be thought-about as 'abstract entities', forming a clone category of 2 cloned files.

## 2.1.4 Crosscutting Structural Clones

Structural clones will crosscut files (or categories, modules etc.), This instance involves 3 PHP files happiness to an online portal module that supports 2 similar crosscutting options. This ends up in 2 structural clones, every consisting of code fragments happiness to at least one of the 2 crosscutting options.

## 2.1.5 Structural Clones Supported Inheritance Hierarchy

The relationship(s) between entities of a structural clone will vary wide. In object-oriented systems, a collection of entities connected by inheritance is accustomed outline a structural clone. we tend to found such a case within the Buffer library (java.nio.*) every consisting of seven Java categories. a lot of data on the structure of the Buffer library, and the way 'feature combinatorics' drawback gave rise to the current structure.

## 3. Detection Structural Clones with Data Processing

This section focuses on our approach to sight some specific sorts of structural clones from all-time low up analysis of similarities employing a data processing technique kind of like the well-known market basked analysis. This analysis builds on the detection of straightforward clones mentioned within the previous chapter. we tend to introduce

associate unvarying approach for the detection of structural clones by moving from the low-level similarities to the upper level similarities. to boost the amount of study, we tend to use abstraction of entities supported clone coverage. the upper level entities that have vital low-level biological research square measure classified along. These teams of entities kind the essential similarity blocks for successive higher level analysis.

## 3.1 Finding Revenant Patterns of Easy Clone Categories

Here we tend to describe the detection of patterns of straightforward clones in a very file - the primary level of structural clones. a similar technique is applied to sight structural clones at alternative levels, as are going to be delineate later. Associate example of this format is shown. once detection easy clone categories in a very system, the info of straightforward clone categories is organized in terms of files depicted by their IDs. the primary information row says that the file with file ID twelve contains 3 clone instances happiness to clone category nine and one instance from every of the clone categories fifteen, 28, 38, and 40. The interpretation is likewise for the opposite information rows. To sight the revenant patterns of straightforward clones in numerous files, we tend to apply the "market basket analysis" technique from the info mining domain. the thought behind this method is to search out the things that square measure sometimes purchased along by totally different customers from a supermarket. These patterns of clone categories can act because the distinctive illustration for a bunch of files, and relying upon its significance in terms of files' coverage, can result in distinguishing teams of extremely similar files. this may be successive level of structural clones. Market basket analysis is completed with "frequent itemset mining (FIM)". The distinction between our drawback and therefore the customary drawback for frequent itemset mining is that in FIM, the things in a very dealings square measure thought-about distinctive, whereas in our information, one file might contain multiple instances of a similar clone category. we are able to normalize the info by removing by doing therefore, we tend to miss out the necessary data, as multiple occurrences of the instances of same clone category in numerous files could be a valid pattern of clones. for instance, 9, 9, 9, fifteen could be a valid clone pattern depicted in all frequent itemsets returns several frequent itemsets that square measure subsets of larger frequent itemsets. the right resolution in our case is to perform "Frequent Closed Itemset Mining" (FCIM), wherever solely those itemsets square measure rumored that don't seem to be subsets of any larger frequent itemset. one in every of the parameters for FCIM is that the support count, which implies the quantity of files that contain the detected pattern of straightforward clones In our case, we've got arduous coded the support to be two, so it'll report a clone pattern, albeit it's gift solely in two files, because it may well be still be vital for maintenance supported its size. The output from FCIM is within the format.

every row represents one frequent clone pattern in conjunction with its support count, indicating the quantity of files containing this clone deals with detection frequent patterns. These clone patterns may also be thought-about as unrestricted gapped clones, wherever any variety of gaps square measure allowed with discretionary size and ordering. a lot of work is needed to isolate clone patterns wherever the gaps square measure little and therefore the clones square measure a lot of cohesive, to possess a lot of pregnant gapped clones. Finding gapped clones during this means conjointly offer the pliability to sight rearranged gapped clones, wherever the cloned elements will occur in discretionary order and will not essentially be organized within the same means.

**3.2 Clump Extremely Cloned Files**

To live file coverage by a clone pattern, we tend to calculate 2 metrics, particularly the File share Coverage (FPC), that indicates the share of a file coated by a clone pattern, and therefore the File Token Coverage (FTC), that tells the quantity of tokens in a very file coated by the clone pattern. These metrics square measure calculated for every file containing the clone pattern. One complication here is that some clones might overlap in a very file, as mentioned earlier, therefore we tend to cannot merely add up the scale of all clones in a very pattern to search out the file coverage. The clump supported these values and alternative parameters may also be made exclusively customizable to suit the wants of the various users. Currently, we tend to let the user specify a minimum FPC and Federal Trade Commission worth to point the importance of a cluster. The cluster are going to be thought-about vital albeit one file has the FPC or Federal Trade Commission worth bigger than threshold values. The expected output is to search out all the numerous clusters that cowl most variety of files and no file is ideally recurrent in 2 clusters. Step one of the rule says that we tend to take away from thought all those clusters wherever no file passes the minimum criteria of FPC and Federal Trade Commission values. These square measure clusters wherever terribly little clones exist between files. In step 2, we tend to type clusters supported their support count. once the support count of 2 or a lot of clusters is same, the clusters square measure sorted supported the utmost FPC worth of the constituent files. Steps three and four prune clusters. These clusters of extremely similar files provide North American nation successive level of structural clones that we tend to decision file clone categories or Federal Communications Commission.

**4. Tool Implementation**

CCFinder implements the structural clone detection techniques given during this paper. CCFinder is written in C++, and it's its own token-based easy clone detector [6]. For frequent closed item sets mining (FCIM), we tend to square measure victimisation the rule from [21] .For manipulation of clones' information, CCFinder makes use of the STL

containers from the quality C++ library. The output from CCFinder is generated within the variety of text files so any visualisation tool developed within the future will simply interface with Clone laborer. For performance analysis, we tend to ran CCFinder on full J2SE one.5 ASCII text file,3 consisting of half-dozen,558 supply files in 370 directories, 625,096 LOC (excluding comments and blank lines), and 70,285 ways, victimisation totally different values of minimum clone size. For forming FCSets and MCSets, a price of fifty tokens is employed for the clump parameter minLen, wherever the Len is measured in terms of tokens. Likewise, for minCover, a price of fifty p.c is employed altogether cases. The tests were run on a Pentium IV machine with three.0 gigahertz processor and 1GB RAM. anytime it took around 2 to a few minutes to run the entire method from finding easy clones to the analysis of files, methods, and directories for structural clones

## 5. Structural Clones In Code

To boost the planning of the inheritance systems, varied re-structuring or refactoring techniques is applied [Opd92] [Fow99]. Analysis of structural clones is useful in locating places wherever high-level duplication is gift, which might be restructured or refactored. For redesigning the system to boost maintainability of the inheritance code, we tend to propose some Structural clone primarily based techniques. an honest place to begin is to investigate the file clone categories (i.e., teams of cloned files). Once selecting file clone categories for refactoring, easy clones or technique clones inside those teams of files is a lot of simply refactored attributable to the context data. it should even be potential to use many little refactoring at the same time, for instance moving along many cloned ways to the parent category, or just ever-changing the inheritance structure to get rid of duplicates. Having solely the information of straightforward clones, risk of constructing such larger changes isn't terribly apparent, and one should go step by step, with the danger of missing the larger image altogether. Analysis may also be done at code fragments, methods, or directories level, looking on however intense the biological research is and the way major reengineering is sensible. the opposite analysis options designed within Clone analyser, like the structural clone configurability and therefore the Diff feature will aid the user within the finer details of refactoring. This planned technique is somewhat general because of the variable objectives; it offers a basic framework for the analysis method.

## 6. Conclusions

We tend to stressed the necessity to check code biological research at a better level. we tend to introduced the idea of structural clone as a continuation configuration of lower-level clones. we tend to given a way for detection structural clones. The method starts by finding easy clones (that is, similar code fragments). more and more higher-level

similarities square measure then found incrementally victimisation data processing technique of finding frequent closed item sets, and clump. we tend to enforced the structural clone detection technique in a very tool referred to as Clone laborer. Whereas Clone laborer may also sight easy clones, its underlying structural clone detection technique will work with the output from any easy clone detector. Structural clone data results in higher program understanding, maintenance, reengineering and use.

**7. References**

1. H.A. Basit and S. Jarzabek, "Detecting Higher-Level Similarity Patterns in Programs," Proc. European code Eng. Conf. and ACM SIGSOFT Symp. Foundations of code Eng., pp. 156-165, Sept. 2005.

2. J.R. Cordy, "Comprehending Reality: sensible Barriers to Industrial Adoption of code Maintenance Automation, "Proc. eleventh IEEE Int'l Workshop Program Comprehension, (keynote paper), pp. 196-206, 2003.

3. A.De Lucia, G. Scanniello, and G. Tortora, "Identifying Clones in Dynamic internet sites victimisation Similarity Thresholds," Proc. Int'l Conf. Enterprise data Systems, pp. 391-396, 2004.

4. J.Y. Gil and I. Maman, "Micro Patterns in Java Code," Proc. twentieth Object bound Programming Systems Languages and Applications, pp. 97-116, 2005.

5. G.Grahne and J. Zhu, "Efficiently victimisation Prefix-Trees in Mining Frequent Itemsets," Proc. 1st IEEE ICDM Workshop Frequent Itemset Mining Implementations, Nov. 2003.

6. J. Han and M. Kamber, "Data Mining: ideas and Techniques". "Morgan dramatist Publishers", 2001.

7. Y. Higo, T. Kamiya, S. Kusumoto, and K. Inoue, "ARIES: Refactoring Support surroundings supported Code Clone Analysis," Proc. Eighth IASTED Int'l Conf. code Eng. and Applications, pp. 222-229, Nov. 2004.

8. S. Jarzabek, "Effective code Maintenance and Evolution: Reused- primarily based Approach". "CRC Press, Taylor and Francis", 2007.

9. S. Jarzabek and S. Li, "Unifying Clones with a Generative Programming Technique: A Case Study," J. code Maintenance and Evolution: analysis and follow, vol. 18, no. 4, pp. 267-292, July 2006.

10. C. Kapser and M.W. Godfrey, "Toward a Taxonomy of Clones in supply Code: A Case Study," Proc. Int'l Workshop Evolution of huge Scale Industrial code Architectures, pp. 67-78, 2003.

11. C. Rich and L.M. Wills, "Recognizing a Program's Design: Ac Graph-Parsing Approach," IEEE code, vol. 7, no. 1, pp. 82-89, Jan. 1990.