

# Software evolution report series 2

---

- Simon Schneider
- Laurance Saess

## TO DO

---

```
1 * Show % of duplicated lines
2 * Show number of clones
3 * Show number of clone classes
4 * Show biggest clone (in lines),
5 * Show biggest clone class
6 * Show example clones
7 * Add raw data table and explain why it was created
```

## Clone detection

---

We use a approach that uses the AST to detect code clones. We use the tactics that is described by [Koschke, 2008]. In pseudo code:

Where:

```
1 * x is the clone type
2 * y is the project location
3 * z is the min number of sub notes
4 * z' is the min number of lines of code per node
```

```
1 type = x
2 a = (for type = 1:100 2:100 3:30)
3
4 ast = loadAstForProject(y);
5 ast <- normalize @ type 2 / type 3
6 astList = getAllNodes(ast)
7 astList <- remove when subitems less than z
8 astList <- remove node when less than z' lines of code
9 astList <- remove with invalid location
10
11 @no duplicate compares
12 @Do not compare when node is a subnode of that node
13 @similarity of nodeA and NodeB > a
14 compare astList astList to nodeA nodeB:
15     return add connection:<nodeA.1,nodeB.1>
16
17 return set nodes:astList
18 return
19
```

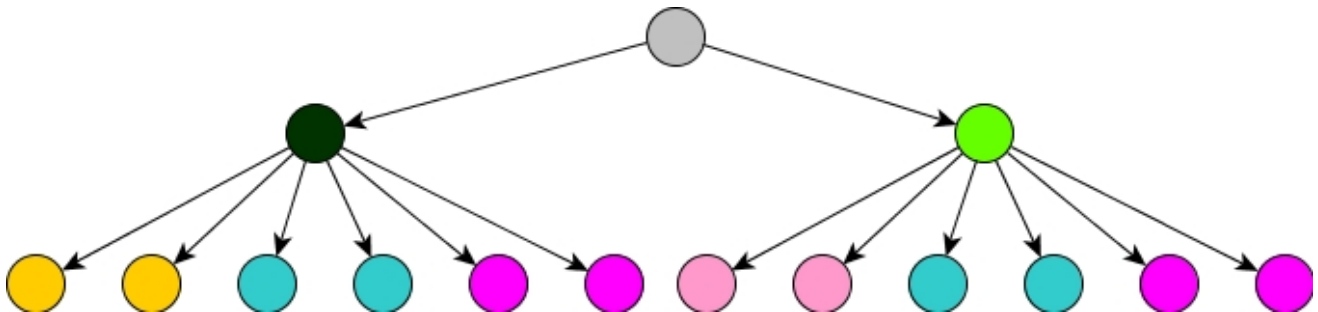
```

20 where:
21     similarity nodea nodeb =
22         nodea.subItems `similar` nodeb.subItems /
23         (nodea.subcount + nodeb.subcount) * 100

```

## Limitations of our AST approach

We cannot compare nodes, or groups of nodes that are on the same level to other nodes or group of nodes. For example, the tree below where nodes with a similar color are the same:



- The two large green trees are not seen as a type one or two clone, this because the sub nodes are different, what is correct.
- The two yellow sub nodes of the first green tree are not seen as duplicates. This because our tool only compares whole nodes.
- The two yellow sub nodes of the first green tree are also not compared with any combination of sequential nodes of the second green sub tree. This because our approach does not support comparing sequential nodes with other nodes in the system. It also does not support comparing these groups to groups under the same parent node.

We did not add support for this because it adds a lot of complexity. The clone detection will also take a lot more times because of the high amount of additional checks that have to be preformed. Large clones are unlikely to be affected because in real-live projects, they normally are not on a single layer in the AST. For this project, and our goal is to find large clones; small codes have a smaller impact to the system.

## Parameters

The real project has the following parameters:

- set[Declaration] ast
- bool normalizeAST
- int minimalNodeGroupSize
- int minimalCodeSize
- real minimalSimilarity

What is a little bit different than the pseudo code. In the section we are going to describe what every parameter is and how it translates to the real project.

## X is the clone type

You can define the clone type in the pseudo code. In the real project you have to translate it like this:

Type 1:

- `normalizeAST = false`
- `minimalSimilarity = 100`

Type 2:

- `normalizeAST = true`
- `minimalSimilarity = 100`

Type 3:

- `normalizeAST = true`
- `minimalSimilarity = 50` <- or any other similarity factor you prefer.

For these settings, `normalizeAST` will remove all information that are type or name related. With this setting `int test` is the same as `float test2`.

For these setting, `minimalSimilarity` will defined a percentage of how much of the node has to be the same to be considered equivalent.

## Y is the project location

The pseudo code will generate an AST based on the location of the project. The clone detection function requires the AST already what is done in the main.

- `ast = createAstsFromEclipseProject(createM3FromEclipseProject(y), true) ;`

## Z is the min number of sub nodes and Z' are the minimum lines of code per node

You can display the AST as an tree. When you compare the nodes, there will be a lot of useless small clones. This parameter can be used to define a minimum size. Nodes are only considered that contain z sub nodes or has minimum z' lines of code.

- `int minimalNodeGroupSize = z`
- `int minimalCodeSize = z'`

# Finding what lines are duplication

---

We use an custom algorithm for detecting the amount of duplicate lines. The algorithm works like this:

1. Get all locations that contain an duplicate
2. Request from the M3 model, all the locations with comments
3. Go through every duplication, and get the lines for every location and look per line:
  1. :If it is a one line comment
  2. If the first character is a multi line start comment
  3. If it ends with an comment close tag

4. If it is an empty line
5. If the line is in a multi-line comment
4. When when of above is true, the line is counted as a non-duplicate
5. When above is false, the line is counted as a duplicate

We wrote the following unit test:

```
1  import javax.annotation.Generated;
2
3  class dupTest {
4      public void testM() {           //Comment
5          int i1 = 1 + 1;             //Comment
6          int i2 = 1 + 1 * 2;         //Comment
7      /*test*/int i3 = 1 + 1 / 4; /*test*/
8  }
9
10     public void testM2() {           //Some test comment
11     /*test*/int i1 = 1 + 1;
12     /*est*/int i2 = 1 + 1 * 2; /*test*/
13         /*
14         * test
15         */
16         int i3 = 1 + 1 / 4;
17     }
18
19     public void testM3() {
20         int i2 = 1 /*test*/+ 1 * 2;
21         //test
22         /*
23
24
25         */
26         int i3 = 1 + 1 / 4;
27         /*
28         int i1 = 1 + 1;
29         int i2 = 1 + 1 * 2;
30         */
31     }
32
33 }
34
```

The file has the following duplicate lines:

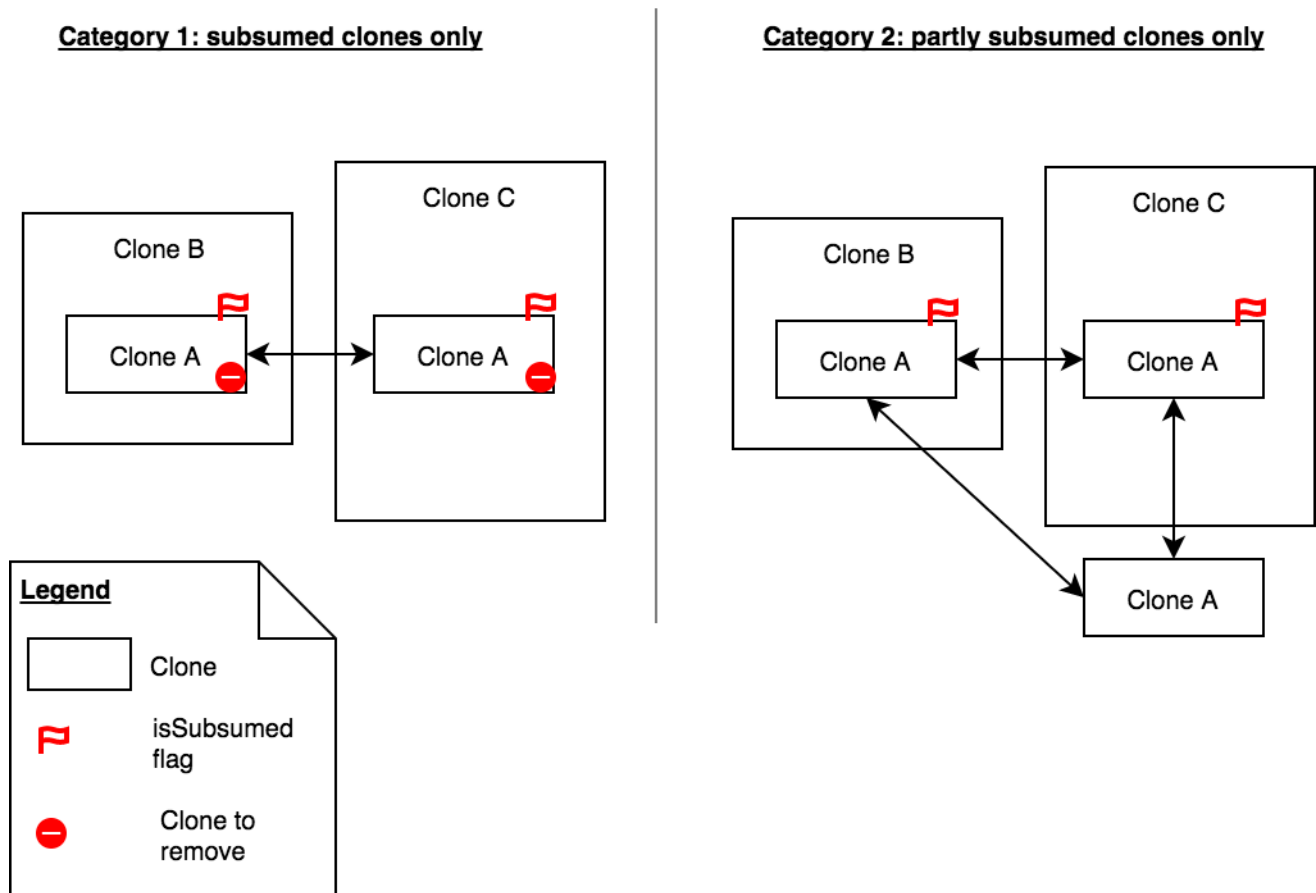
```
1  [4, 5, 6, 7, 8, 10, 11, 12, 16, 17, 20, 26]
```

What is correct.

## Clone classes

The diagram below shows our approach to delete unnecessary clones, which are subsumed by other clones.

**Figure: Detecting clone classes**



At first, we flagged every clone that was included in another clone. For this, we can use the underlying location of a cloned node.

```
1 public bool containsLargerLocation(NodeIdLocation n1A, list[NodeIdLocation] n1s){
2     return any(NodeIdLocation n1B <- n1s, n1A.id != n1B.id && n1B.l >= n1A.l);
3 }
```

This code above can tell us if a node (`n1A`) is located inside a clone of the node list (`n1s`). In this first approach, we simply flagged and deleted the clones from the clone graph.

We soon recognized another type of clones (category 2), in which included clones are not deleted because they appear outside of other clones. In order to build the right clone classes, we use the following algorithm:

1. Flag every clone that is fully included in another clone.
2. Build a graph of clones
  1. Use the transitive closure for clone type 1 and 2. Type3 clones must not be transitive.
3. Remove a class of clones from the output if all nodes are flagged as deleted.
4. Output the clones as a graph that builds a cluster for every clone class (please see visualization below)

## Visualization

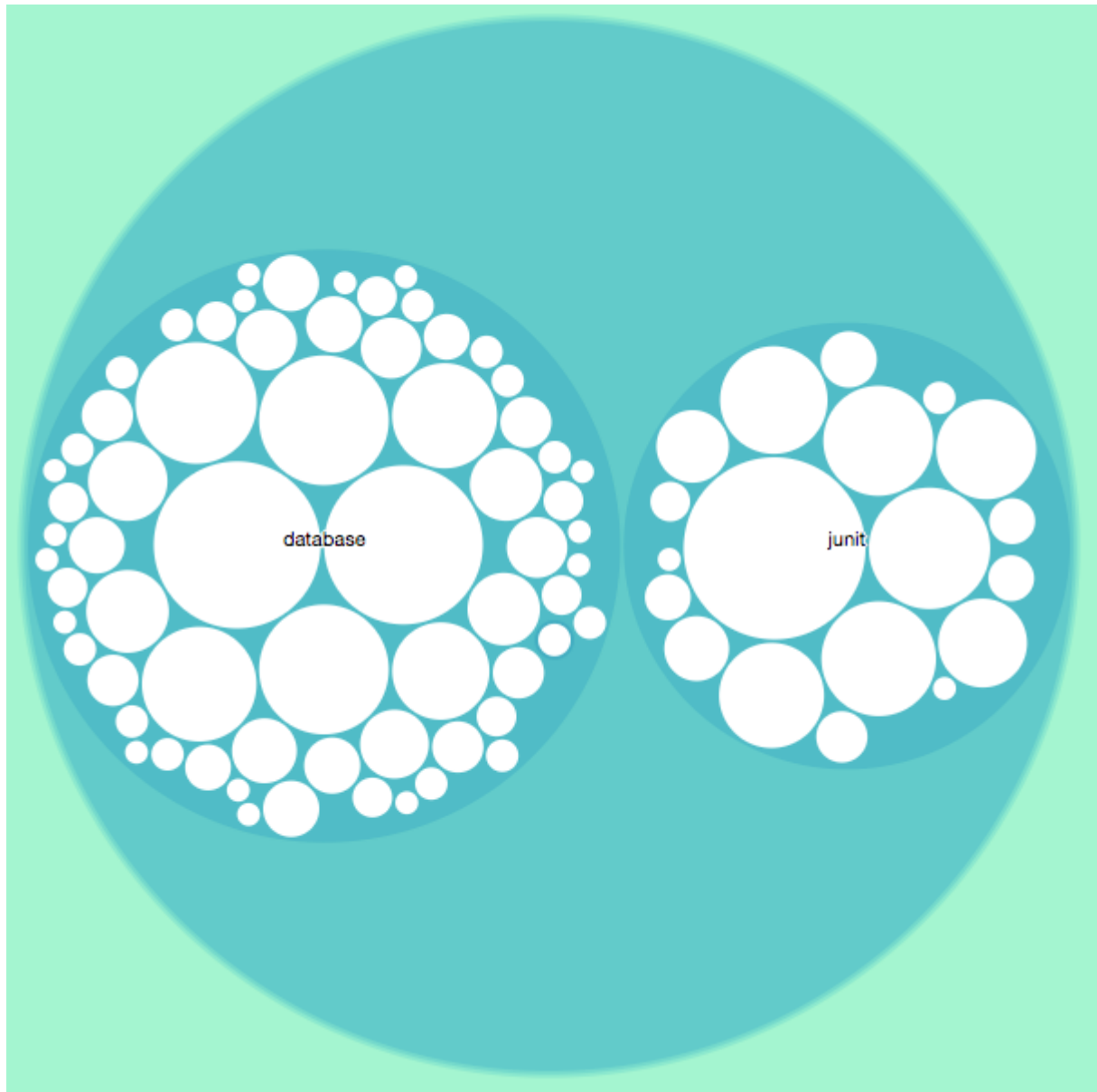
## Enclosure diagram

---

In his book, Adam Tornhill describes the methods and visualizations that criminologists use to find murderers and describes how software engineers can learn from this, by using similar methods to detect code smells or "dirty" developers. [Tornhill, 2015]

One example of visualizations that criminologists supposedly use maps in which they highlight hotspots in order to predict future crimes.

**Figure: Clone enclosure diagram to show clone overview**



The visualization above tries to imitate this approach by creating a zoomable map of code clone criminals.

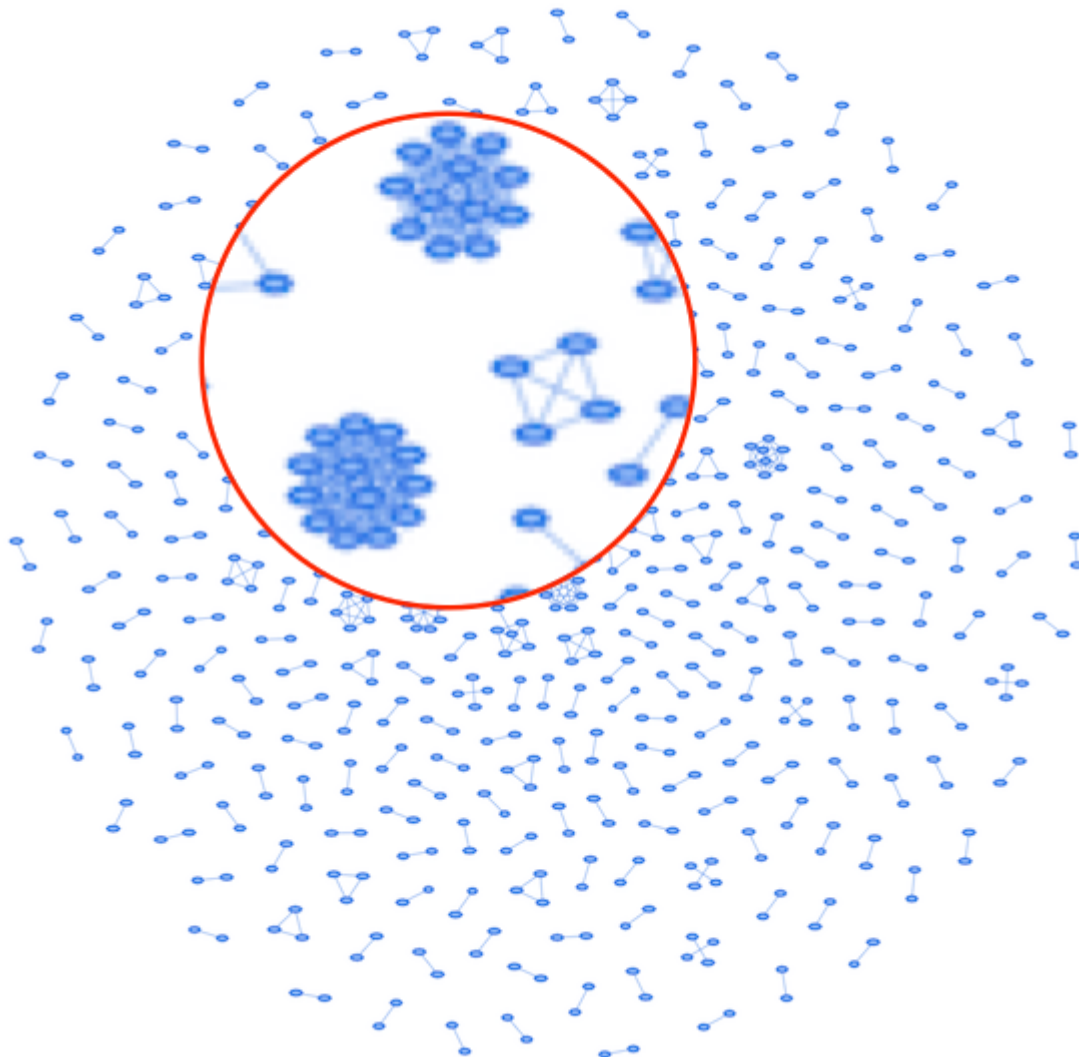
We used an enclosure diagram in order to handle large software systems. This diagram is based on a geometric layout algorithm called circle packing. Each circle represents a module or code class of the system. The more duplicate lines a module has the larger the circle. We can immediately see the files with the largest clones and in which folders they are located.

## Clone graphs

---

The enclosure diagram is good to get a quick overview over the project, but it does not tell anything about how clones are connected inside the system.

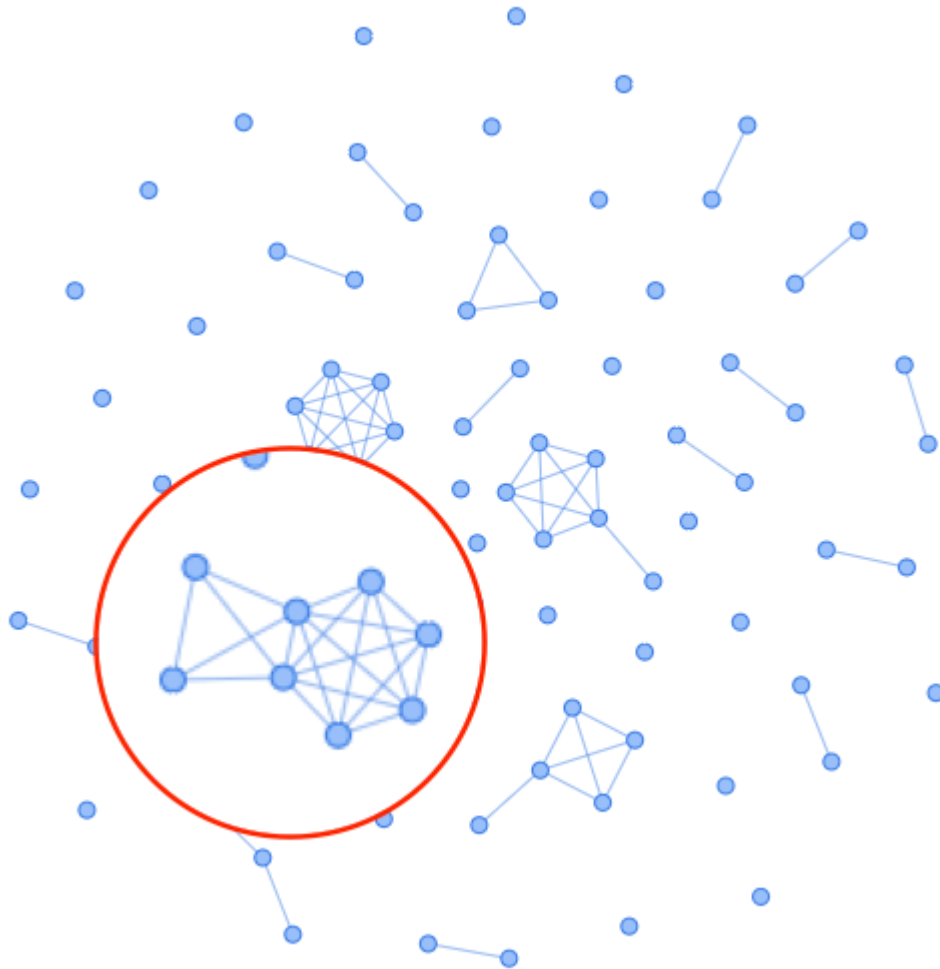
**Figure: Graph of clone classes (approach 1)**



To show the maintainer how the clone classes in his project are structured we created the interactive diagram above.

In the first approach, we used the clones as nodes and their connections as edges. If there are two clones of type A in one file, they will be displayed as two connected nodes. This view nicely shows that clone classes can be seen as clusters that are not interconnected. A behavior that we used to detect them (see the pseudo algorithm in "clone classes" section).

**Figure: Graph of clone files (approach 2)**



The second approach is more useful for a maintainer. In this view, the files are represented as nodes and they are connected by an edge as soon as one duplicate is found in both files. A highlighted area shows that these clusters are not fully interconnected, unlike in the first approach where the clones were fully transitive. If `class A` and `class B` share a clone, and `class B` and `class C` share a clone, `class A` and `class C` are not connected (not transitive).

What a maintainer can learn from this view:

- There are many files that have clones in themselves - which should be an easy fix
- Most of the clusters are fully interconnected, which shows that there is some behavior that all of these files could have in common.
- Removing the largest 4 clone classes could reduce the number of files that contain clones drastically.

## Maintainer requirements

1: Maintainers have to have an understanding of a program. There are multiple ways of getting an understanding at a global level of the program. In the paper [Storey, Fracchia, Müller, 1999] these are Macro-strategies. One strategy is called the "As-needed macro-strategies" strategy [Storey, Fracchia, Müller, 1999]. You only take a look at a part of the code when you needed it. With this approach, code clones give a very negative impact. When you want to make a change, you first have to find where the duplicates are,



otherwise you can introduce bugs into the system. Our goal is to visualization of clone classes. We are going to resolve this by adding support for clicking so that you can see where else the same code is used.

2: One issue in maintenance can be that duplicate in required to extent current functionality. When a maintainer wants to abstract these parts, he has to know on what locations these abstractions can be introduced. When you can find clones, and view how to clones look like, you can faster spot places where you can introduce abstraction. This is why we are making it possible to show the code per duplicate. Clicking in the clone will result in a overview of the code.

3: When a maintainer wants to improve the code quality, it is useful to know in what part of the project the most duplicates are. This way, you know on what part of the project the most progress can be booked. It can be use as a guide for reducing complexity. To solve this, we added an overview of the largest code classes.

---

## Literature

---

[Tornhill, 2015] Adam Tornhill. (2015). Your Code as a Crime Scene.

<https://doi.org/10.1017/CBO9781107415324.004>

[Koschke, 2008] R. Koschke. (2008). Identifying and Removing Software Clones.

[Storey, Fracchia, Müller, 1999] Storey, M. A., Fracchia, F. D., & Müller, H. A. (1999). Cognitive design elements to support the construction of a mental model during software exploration. *Journal of Systems and Software*, 44(3), 171-185.