

Methods and Tools for the Analysis of Legacy Software Systems

Department of Computers and Information Technology
2021

Ph.D. student: Stana Adelina Diana

Scientific supervisor: prof.dr.ing Cretu Vladimir-Ioan

Contents

1	Extracting software dependencies	4
1.1	Tool for measuring software dependencies	4
1.2	Extracting structural dependencies	5
1.3	Extracting logical dependencies	6
2	Filtering extracted logical dependencies	8
2.1	Data set used	8
2.2	Filtering based on the size of commit transactions	10
2.3	Filtering based on the number of occurrences	11
2.4	Overlaps between structural and logical dependencies	17
3	Usage of the extracted dependencies	21
3.1	Data set used	21
3.2	Identifying key classes using logical dependencies	24
3.2.1	Definition and previous work	24
3.2.2	Metrics for key classes detection	25
3.2.3	Metrics for results evaluation	27
3.2.4	Previous measurements	28
3.2.5	Measurements using logical dependencies	29

3.3	Comparison of the extracted data with fan-in and fan-out metric . . .	32
-----	---	----

Chapter 1

Extracting software dependencies

1.1 Tool for measuring software dependencies

In order to build structural and logical dependencies we have developed a tool that takes as input the source code repository and builds the required software dependencies [33]. The workflow can be delimited by three major steps as it follows (Figure 1-1):

Step 1: *Extracting structural dependencies.*

Step 2: *Extracting logical dependencies.*

Step 3: *Processing the information extracted.*

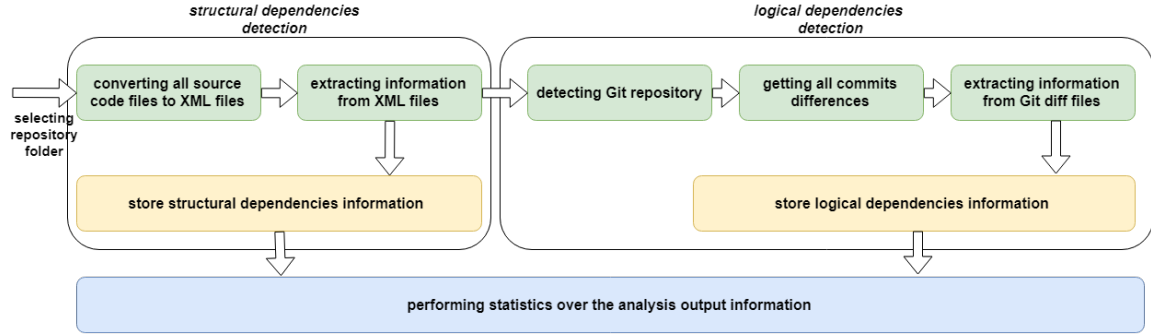


Figure 1-1: Processing phases

1.2 Extracting structural dependencies

A dependency is created by two elements that are in a relationship and indicates that an element of the relationship, in some manner, depends on the other element of the relationship [6], [10].

Structural dependencies can be found by analyzing the source code [30], [8], [5]. A structural dependency between two classes A and B is given by the fact that A statically depends on B, meaning that A cannot be compiled without knowing about B. In object oriented system, this dependency can be given by many types of relationships between the two classes: A extends B, A implements B, A has attributes of type B, A has methods which have type B in their signature, A uses local variables of type B, A calls methods of B.

We use an external tool called srcML [11], [12] to convert all source code files from the current release into XML files. All the information about classes, methods, calls to other classes are afterwards extracted by our tool parsing the XML files and building a dependencies data structure. We have chosen to rely on srcML as a preprocessing tool because it reduces a significant number of syntactic differences

from different programming languages and can make easier the parsing of source code written in different programming languages such as Java, C++ and C#.

1.3 Extracting logical dependencies

Logical dependencies (a.k.a logical coupling) can be found by software history analysis and can reveal relationships that are not always present in the source code (structural dependencies).

Software engineering practice has shown that sometimes modules which do not present structural dependencies still appear to be related. Co-evolution represents the phenomenon when one component changes in response to a change in another component [38], [9]. Those changes can be found in the software history maintained by the versioning system. Gall [18], [19] identified as logical coupling between two modules the fact that these modules *repeatedly* change together during the historical evolution of the software system [3].

The versioning system contains the long-term change history of every file. Each project change made by an individual at a certain point of time is contained into a commit [13]. All the commits are stored in the versioning system chronologically and each commit has a parent. The parent commit is the baseline from which development began, the only exception to this rule is the first commit which has no parent. We will take into consideration only *commits that have a parent* since the first commit can include source code files that are already in development (migration from one versioning system to another) and this can introduce redundant logical links [1].

The tool looks through the main branch of the project and gets all the existing commits. For each commit a diff against the parent will be made and stored. Here

we have the option to ignore commits that contain more files than a threshold value for commit size. Also, we have the option to check whether the differences are in actual code or if they affect only parts of source files that are only comments. Finally after all the difference files are stored, all the files are parsed and logical dependencies are build. For a group of files that are committed together, logical dependencies are added between all pairs formed by members of the group. Adding a logical dependency increases an occurrence counter for the logical link.

Chapter 2

Filtering extracted logical dependencies

2.1 Data set used

We have analyzed a set of open-source projects found on GitHub¹ [20] in order to extract the structural and logical dependencies between classes. Table 2.1 enumerates all the systems studied. The 1st column assigns the projects IDs; 2nd column shows the project name; 3rd column shows the number of entities(classes and interfaces) extracted; 4th column shows the number of most recent commits analyzed from the active branch of each project and the 5th shows the language in which the project was developed.

¹<http://github.com/>

Table 2.1: Summary of open source projects studied.

ID	Project	Nr. of entites	Nr. of commits	Type
1	bluecove	586	894	java
2	aima-java	987	818	java
3	powermock	1084	893	java
4	restfb	783	1188	java
5	rxjava	2673	2468	java
6	metro-jax-ws	1103	2222	java
7	mockito	1409	1572	java
8	grizzly	1592	3122	java
9	shipkit	242	1483	java
10	OpenClinica	1653	3749	java
11	robolectric	2050	5029	java
12	aeron	541	5101	java
13	antlr4	1381	3449	java
14	mcidasv	805	3668	java
15	ShareX	919	2505	csharp
16	aspnetboilerplate	2353	1615	csharp
17	orleans	3485	3353	csharp
18	cli	767	2397	csharp
19	cake	2250	1853	csharp
20	Avalonia	1677	2445	csharp
21	EntityFramework	7107	2443	csharp
22	jellyfin	2179	4065	csharp
23	PowerShell	861	2033	csharp
24	WeiXinMPSDK	2029	2723	csharp
25	ArchiSteamFarm	117	2181	csharp
26	VisualStudio	1016	4417	csharp
27	CppSharp	259	3882	csharp

2.2 Filtering based on the size of commit transactions

A big commit transaction can indicate that a merge with another branch or that a renaming has been made. In this case, a series of irrelevant logical dependencies can be introduced since not all the files are updated in the same time for a development reason. Different works have chosen fixed threshold values for the maximum number of files accepted in a commit. Cappiluppi and Ajienka, in their works [1], [2] only take into consideration commits with less than 10 source code files changed in building the logical dependencies.

The research of Beck et al [4] only takes in consideration transactions with up to 25 files. The research [24] provided also a quantitative analysis of the number of files per revision; Based on the analysis of 40,518 revisions, the mean value obtained for the number of files in a revision is 6 files. However, standard deviation value shows that the dispersion is high.

We analyzed the overall transaction size trend for 27 open-source cpp and java systems. The results are presented in Figure 2-1, based on them we can say that 90% of the total commit transactions made are with less than 10 source code files changed. This percent allows us to say that setting a threshold of 10 files for the maximum size of the commit transactions will not affect so much the total number of commit transactions from the systems since it will still remain 90% of the commit transactions from where we can extract logical dependencies [33].

As we can see in Figure 2-2 even though only 5% of the commit transactions have more than 20 files changed ($20 < cs < inf$) they generate in average 80% of the total amount of logical dependencies extracted from the systems. The high number

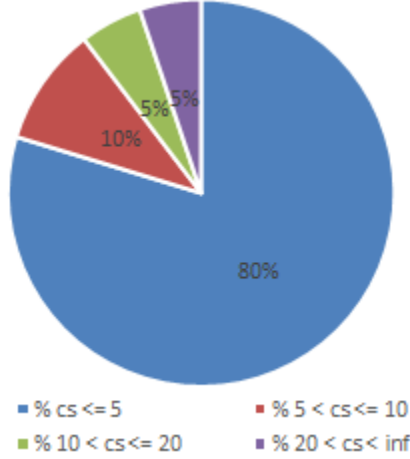


Figure 2-1: Commit transaction size(cs) trend in percentages

of logical dependencies extracted from such a small number of commit transactions is caused by big commit transactions. One single big commit transaction can lead to a large amount of logical dependencies. For example in RxJava we have a very few commit transactions with 1030 source code files, this means that those files can generate ${}^nC_k = \frac{n!}{k!(n-k)!} = \frac{1030!}{2!(1028)!} = 529935$ logical dependencies. By setting a threshold on the commit transaction size we can avoid the introduction of those logical dependencies into the system.

So filtering 10% of the total amount of commit transactions can indeed lead to a significant decrease of the amount of logical dependencies and that is why we choose the value of 10 files as our fixed threshold for the maximum size of a commit transaction [33].

2.3 Filtering based on the number of occurrences

One occurrence of a co-change between two software entities can be a valid logical dependency, but can also be a coincidence. Taking into consideration only co-changes

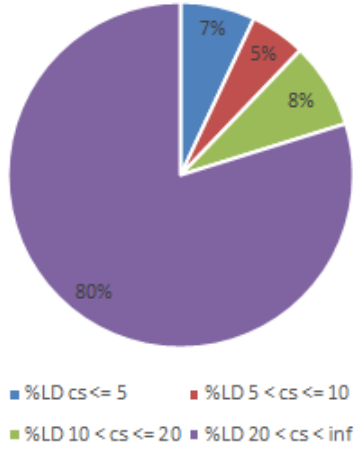


Figure 2-2: Percentages of LD extracted from each commit transaction size(cs) group

with multiple occurrences as valid dependencies can lead to more accurate logical dependencies and more accurate results. On the other hand, if the project studied has a relatively small amount of commits, the probability to find multiple updates of the same classes in the same time can be small, so filtering after the number of occurrences can lead to filtering all the logical dependencies extracted. Giving the fact that we will study multiple projects of different sizes and number of commits, we will analyze also the impact of this filtering on different projects.

We have performed a series of analysis on the test systems, incrementing the threshold value *occ* from 1 to 4. In each of the cases the extracted logical dependencies from commit transaction with less or equal to 10 changed source code files were also filtered by the minimum number of occurrences established and all the logical dependencies that did not exceeded the minimum number of occurrences were discarded.

The results of the analysis are presented in Table 2.2 as percentages of logical dependencies (LD) that are also structural dependencies and Table 2.3 as ratio of

the number of logical dependencies (LD) to the number of structural dependencies (SD).

Table 2.2: Percentage of LD that are also SD

ID	$occ \geq 1$	$occ \geq 2$	$occ \geq 3$	$occ \geq 4$
1	7,13	7,77	7,99	19,71
2	19,54	25,76	29,55	32,16
3	6,66	8,58	11,82	14,87
4	1,16	1,17	0,91	0,80
5	3,99	3,96	7,75	7,49
6	13,92	20,16	22,91	22,77
7	8,38	9,28	14,93	14,58
8	6,70	9,73	14,20	15,60
9	16,98	23,34	29,22	32,89
10	8,94	9,15	11,05	10,59
11	4,99	6,92	8,88	11,08
12	13,19	17,15	18,60	19,57
13	2,43	5,59	8,33	8,21
14	13,27	18,88	19,02	19,28
15	12,90	21,95	25,51	27,01
16	13,33	17,34	18,53	16,24
17	6,09	6,18	6,41	6,44
18	9,73	10,60	14,27	18,80
19	10,26	13,54	13,64	12,60
20	12,83	18,36	21,00	25,72
21	2,86	4,65	5,70	4,98
22	5,20	6,56	8,18	8,90
23	8,23	13,64	17,04	17,65
24	6,77	10,89	14,47	16,05
25	9,85	10,14	11,65	11,33
26	8,65	10,79	12,78	14,34
27	7,04	8,78	9,87	10,08
Avg	8,93	11,88	14,23	15,55

Table 2.3: Ratio of number of LD to number of SD

ID	$occ \geq 1$	$occ \geq 2$	$occ \geq 3$	$occ \geq 4$
1	4,13	1,94	1,23	0,26
2	0,81	0,33	0,16	0,10
3	5,12	1,93	0,78	0,38
4	53,36	42,00	38,31	36,30
5	4,27	2,90	0,88	0,72
6	1,07	0,46	0,30	0,23
7	4,09	2,38	0,99	0,73
8	4,06	1,57	0,76	0,49
9	3,64	2,03	1,14	0,77
10	1,41	1,01	0,47	0,34
11	7,91	4,47	2,93	2,03
12	3,92	2,15	1,47	1,07
13	10,15	3,18	1,22	1,03
14	3,07	1,53	1,16	0,97
15	2,34	0,84	0,48	0,33
16	1,21	0,47	0,26	0,19
17	2,99	1,83	1,11	0,84
18	2,26	1,37	0,67	0,40
19	2,32	1,38	0,76	0,67
20	1,24	0,58	0,35	0,18
21	5,33	2,12	1,27	1,05
22	3,38	1,88	0,99	0,74
23	3,62	1,22	0,76	0,37
24	2,57	1,22	0,67	0,46
25	7,47	5,36 ₁₅	4,16	3,73
26	4,03	2,16	1,50	1,15
27	7,46	4,26	2,99	2,43
Avg	5,67	3,43	2,51	2,15

Based on Table 2.2 we can say that only a small percentage of the extracted logical dependencies are also structural dependencies. This is consistent with the findings of related works [1], [2]. The percentage of LD which are also SD increases with the minimum number of occurrences because the number of logical dependencies from the systems decreases with the minimum number of occurrences. We calculate the overlapping between logical and structural dependencies not only because we want to get an idea of how many structural dependencies are reflected in the versioning system through logical dependencies but also because we want to eliminate logical dependencies that are also structural dependencies since they don't bring any new information to the systems.

We stopped the minimum occurrences threshold to 4 because we observed that for systems with ID 2, 6, 10 and 16 from Table 2.3 the ratio number is lower than 1 which means that the number of SD is higher than the number of LD. On the other hand for systems with ID 4, 11, 25, 27 the threshold of 4 for minimum number of occurrences does not change the discrepancy between the number of logical and structural dependencies. If we try to go higher with the occurrences threshold we will risk to filter all the existing logical dependencies for some of the systems. So, filtering with a threshold of 4 for the minimum number of occurrences will indeed filter the logical dependencies but for some of the systems the remaining number of logical dependencies will still be significantly higher compared to the number of structural dependencies.

2.4 Overlaps between structural and logical dependencies

A logical dependency can be also a structural dependency and vice-versa, so studying the overlapping between logical and structural dependencies while filtering is important since the intention is to introduce those logical dependencies among with structural dependencies in architectural reconstruction systems. Current studies have shown a relatively small percentage of overlapping between them with and without any kind of filtering [1]. This means that a lot of non related entities update together in the versioning system, the goal here is to establish the factors that determine such a small percentage of overlapping [32].

In the main series of experiments, for each system, we extracted the structural dependencies and the logical dependencies and determined the overlap between the two dependencies sets, in various experimental conditions.

One variable experimental condition is whether changes located in comments contribute towards logical dependencies. This condition distinguishes between two different cases:

- with comments: a change in source code files is counted towards a logical dependency, even if the change is inside comments in all files
- without comments: commits that changed source code files only by editing comments are ignored as logical dependencies

In all cases, we varied the following threshold values:

- commit size (*cs*): the maximum size of commit transactions which are accepted to generate logical dependencies. The values for this threshold were 5, 10, 20

and no threshold (infinity).

- number of occurrences (*occ*): the minimum number of repeated occurrences for a co-change to be counted as logical dependency. The values for this threshold were 1, 2, 3 and 4.

The six tables below present the synthesis of our experiments. We have computed the following values:

- the mean ratio of the number of logical dependencies (LD) to the number of structural dependencies (SD)
- the mean percentage of structural dependencies that are also logical dependencies (calculated from the number of overlaps divided to the number of structural dependencies)
- the mean percentage of logical dependencies that are also structural dependencies (calculated from the number of overlaps divided to the number of logical dependencies)

In all the six tables, 2.4, 2.5, 2.6, 2.7, 2.8, 2.9 we have on columns the values used for the commit size *cs*, while on rows we have the values for the number of occurrences threshold *occ*. The tables contain median values obtained for experiments done under all combinations of the two threshold values, on all test systems. In all tables, the upper right corner corresponds to the most relaxed filtering conditions, while the lower left corner corresponds to the most restrictive filtering conditions.

In order to assess the influence of comments, we compare pairwise Tables 2.4 and 2.5, Tables 2.6 and 2.7 and Tables 2.8 and 2.9. We observe that, although there are some differences between pairs of measurements done in similar conditions with and without comments, the differences are not significant.

Table 2.4: Ratio of number of LD to number of SD, case with comments

	$cs \leq 5$	$cs \leq 10$	$cs \leq 20$	$cs < \infty$
$occ \geq 1$	3,39	5,67	9,00	80,31
$occ \geq 2$	2,24	3,47	5,02	60,14
$occ \geq 3$	1,04	2,53	3,52	44,68
$occ \geq 4$	0,90	2,16	2,88	33,47

Table 2.5: Ratio of number of LD to number of SD, case without comments

	$cs \leq 5$	$cs \leq 10$	$cs \leq 20$	$cs < \infty$
$occ \geq 1$	3,24	5,33	7,90	67,16
$occ \geq 2$	1,35	3,27	4,72	47,39
$occ \geq 3$	1,00	1,67	2,49	32,39
$occ \geq 4$	0,43	1,26	1,93	22,15

Table 2.6: Percentage of SD that are also LD, case with comments

	$cs \leq 5$	$cs \leq 10$	$cs \leq 20$	$cs < \infty$
$occ \geq 1$	19,75	29,86	39,29	76,59
$occ \geq 2$	12,50	20,20	27,68	66,11
$occ \geq 3$	8,49	14,22	19,94	55,99
$occ \geq 4$	6,58	10,95	15,76	47,12

Table 2.7: Percentage of SD that are also LD, case without comments

	$cs \leq 5$	$cs \leq 10$	$cs \leq 20$	$cs < \infty$
$occ \geq 1$	18,88	28,47	37,44	71,12
$occ \geq 2$	11,87	19,03	25,93	59,58
$occ \geq 3$	8,00	13,09	18,15	48,65
$occ \geq 4$	5,85	9,94	14,27	39,07

Table 2.8: Percentage of LD that are also SD, case with comments

	$cs \leq 5$	$cs \leq 10$	$cs \leq 20$	$cs < \infty$
$occ \geq 1$	12,02	8,86	6,72	1,79
$occ \geq 2$	15,05	11,71	9,38	2,21
$occ \geq 3$	17,45	13,97	11,57	2,86
$occ \geq 4$	18,96	15,28	12,94	3,67

Table 2.9: Percentage of LD that are also SD, case without comments

	$cs \leq 5$	$cs \leq 10$	$cs \leq 20$	$cs < \infty$
$occ \geq 1$	12,05	9,02	6,98	1,93
$occ \geq 2$	15,08	12,03	9,66	2,42
$occ \geq 3$	17,78	14,37	12,24	3,28
$occ \geq 4$	19,22	15,59	13,30	4,21

On the other hand, the overlap between structural and logical dependencies is given by the number of pairs of classes that have both structural and logical dependencies. We evaluate this overlap as a percentage relative to the number of structural dependencies in Tables 2.6 and 2.7, respectively as a percentage relative to the number of logical dependencies in Tables 2.8 and 2.9.

A first observation from Tables 2.6 and 2.7 is that not all pairs of classes with structural dependencies co-change. The biggest value for the percentage of structural dependencies that are also logical dependencies is 76.5% obtained in the case when no filterings are done.

From Tables 2.8 and 2.9 we notice that the percentage of logical dependencies which are also structural is always low to very low. This means that most co-changes are recorded between classes that have no structural dependencies to each other [32].

Chapter 3

Usage of the extracted dependencies

3.1 Data set used

To extract the key classes based on logical dependencies, we took the same set of data used in another research involving key class detection. The research of I. Sora et al [26] takes into consideration structural public dependencies that are extracted using static analysis techniques and was performed on the object-oriented systems presented in table 3.1.

The requirements for a system to qualify as suited for investigations using logical dependencies are: has to be on GitHub, has to have release tags to identify the version, and also has to have an increased number of commits. From the total of 14 object-oriented systems listed in the paper [26], 13 of them have repositories in Github 3.2. And from the found repositories we identified only 6 repositories that have the same release tag as the specified version from table 3.1. It is important

to identify the correct release tag for each repository to limit the commits further analyzed by date. Only commits that were made until the specified release are considered and analyzed. The commits number found on the remaining 6 repositories varies from 19108 commits for Tomcat Catalina to 149 commits for JHotDraw. In order to have more accurate results, we need a significant number of commits, so we reached the conclusion that only 3 systems can be used for key classes detection using logical dependencies: Apache Ant, Hibernate, and Tomcat Catalina. From all the systems mentioned in table 3.1 Apache Ant is the most used and analyzed in other works [32], [16], [39], [21].

Table 3.1: Analyzed software systems in previous research paper.

ID	System	Description	Version
S1	Apache Ant	Java library and command line tool that drive the build processes as targets and extension points depending upon each other	1.6.1
S2	Argo UML	UML modelling tool with support for all UML diagrams.	0.9.5
S3	GWT Portlets	Open source web framework for building GWT (Google Web Toolkit) Applications.	0.9.5 beta
S4	Hibernate	Persistence framework for Java.	5.2.12
S5	javacient	Java distributed application for playing with robots	2.0.0
S6	jEdit	Java mature text editor for programmers.	5.1.0
S7	JGAP	Genetic Algorithms and Genetic Programming Java library.	3.6.3
S8	JHotDraw	JHotDraw is a two-dimensional graphics framework for structured drawing editors that is written in Java.	6.0b.1
S9	JMeter	JMeter is a Java application designed to load test functional behavior and measure performance	2.0.1
S10	Log4j	Logging Service	2.10.0
S11	Mars	The Mars Simulation Project is a Java project that models and simulates human settlements on Mars planet	3.06.0
S12	Maze	The Maze-solver project simulates an artificial intelligence algorithm on a maze	1.0.0
S13	Neuroph	Neuroph is a Java neural network framework.	2.2.0
S14	Tomcat Catalina	The Apache Tomcat project is an open-source implementation of JavaServlet and JavaServerPages technologies	9.0.4
S15	Wro4J	The Wro4J is a web resource (JS and CSS) optimizer for Java.	1.6.3

Table 3.2: Found systems and versions of the systems in GitHub.

ID	System	Version	Release Tag name	Commits number
S1	Apache Ant	1.6.1	rel/1.6.1	6713
S2	Argo UML	0.9.5	not found	0
S3	GWT Portlets	0.9.5 beta	not found	0
S4	Hibernate	5.2.12	5.2.12	6733
S5	javaclient	2.0.0	not found	0
S6	jEdit	5.1.0	not found	0
S7	JGAP	3.6.3	not found	0
S8	JHotDraw	6.0b.1	not found	149
S9	JMeter	2.0.1	v2_1_1	2506
S10	Log4j	2.10.0	v1_2_10-recalled	634
S11	Mars	3.06.0	not found	0
S12	Maze	1.0.0	not found	0
S13	Neuroph	2.2.0	not found	0
S14	Tomcat Catalina	9.0.4	9.0.4	19108
S15	Wro4J	1.6.3	v1.6.3	2871

3.2 Identifying key classes using logical dependencies

3.2.1 Definition and previous work

Zaidman et al [40] were the first to introduce the concept of key classes and it refers to classes that can be found in documents written to provide an architectural overview of the system or an introduction to the system structure. Tahvildari and Kontogiannis have a more detailed definition regarding key classes concept: “Usually, the most important concepts of a system are implemented by very few key classes

which can be characterized by the specific properties. These classes, which we refer to as key classes, manage many other classes or use them in order to implement their functionality. The key classes are tightly coupled with other parts of the system. Additionally, they tend to be rather complex, since they implement much of the legacy system’s functionality” [35]. Also, other researchers use a similar concept as the one defined by Zaidman but under different terms like important classes [22] or central software classes [34].

In previous works, the approach for finding key classes is based on ranking the classes with a page ranking algorithm [14], [25], [26], [31] . The page ranking algorithm is a customization of PageRank, the algorithm used to rank web pages [28]. The PageRank algorithm works based on a recommendation system. If one node has a connection with another node, then it recommends the second node. In previous works, connections are established based on structural dependencies extracted from static code analysis. If A has a structural dependency with B, then A recommends B, and also B recommends A.

3.2.2 Metrics for key classes detection

In order to identify the key classes of an object-oriented system, we have to determine what metrics can be used in order to get a good overview of the system and its most important classes [15], [40], [29] . The metrics used in previous research can be grouped into the following categories:

- class size metrics: number of fields (NoF), number of methods (NoM), global size ($\text{Size} = \text{NoF} + \text{NoM}$).
- class connection metrics, any structural dependency between two classes:

- CONN-IN, the number of distinct classes that use a class;
 - CONN-OUT, the total number of distinct classes that are used by a class;
 - CONN-TOTAL, the total number of distinct classes that a class uses or are used by a class (CONN-IN + CONN-OUT).
 - CONN-IN-W, the total weight of distinct classes that use a class.
 - CONN-OUT-W, the total weight of distinct classes that are used by a class.
 - CONN-TOTAL-W, the total weight of all connections of the class (CONN-IN-W + CONN-OUT-W) [26].
- class pagerank values, previous research use pagerank values computed on both directed and undirected, weighted and unweighted graphs:
 - PR - value computed on the directed and unweighted graph;
 - PR-W - value computed on the directed and weighted graph;
 - PR-U - value computed on the undirected and unweighted graph;
 - PR-U-W - value computed on the undirected and weighted graph;
 - PR-U2-W - value computed on the weighted graph with back-recommendations [14], [25], [26], [31].

Because the extracted logical dependencies from the systems are undirected, from the mentioned metrics, we can use the following ones: CONN-TOTAL, CONN-TOTAL-W, PR-U, PR-U-W and PR-U2-W.

3.2.3 Metrics for results evaluation

A classification model is a mapping between expected results and predicted results [17], [7]. Both results can be labeled as positive or negative, which leads us to the confusion matrix from figure 3-1.

Expected Result \ Predicted Result	Positive	Negative
Positive	<i>True Positive</i>	<i>False Positive</i>
Negative	<i>False Negative</i>	<i>True Negative</i>

Figure 3-1: Confusion matrix

The confusion matrix has the following outcomes:

- *true positive*, if the expected result is positive and the predicted result is also positive.
- *false positive*, if the expected result is positive but the predicted result is negative.
- *false negative*, if the expected result is negative but the predicted result is positive.
- *true negative*, if the expected result is negative and the predicted result is also negative.

The true positive rate of a classifier is calculated as the division between the number of true positive results identified and all the positive results identified:

$$\text{True positive rate}(TPR) = \frac{TP}{TP + FN}$$

The false positive rate of a classifier is calculated as the division between the number of false positive results identified and all the negative results identified:

$$\text{False positive rate}(FPR) = \frac{FP}{FP + TN}$$

To calculate the performance of a classification model, the Receiver Operating Characteristic (ROC) graph can be used. The ROC graph is a two-dimensional graph that has on the X-axis plotted the false positive rate and on the Y-axis the true positive rate. By plotting the true positive rate and the false positive rate at thresholds that vary between a minimum and a maximum possible value we obtain the ROC curve. The area under the ROC curve is called Area Under the Curve (AUC).

In multiple related works, the ROC-AUC metric has been used to evaluate the results for finding key classes of software systems [27], [26], [36], [37].

For a classifier to be considered good, its ROC-AUC metric value should be as close to 1 as possible, when the value is 1 then the classifier is considered to be perfect.

3.2.4 Previous measurements

In the research of I. Sora et al [26] is used a tool that takes as an input the source code of the system and applies a ranking strategy to rank the classes according to

their importance. To differentiate the important classes from the rest of the classes, a TOP threshold for the top classes found is set. The threshold can vary between 20 and 30 classes.

The expected results from the research are based on classes labeled as important classes in the system documentation. The true positives (TP) are the classes found in the reference solution and also in the top TOP ranked classes. False positives (FP) are the classes that are not in the reference solution but are in the TOP ranked classes. True Negatives (TN) are classes that are found neither in the reference solution nor in the TOP ranked classes. False Negatives (FN) are classes that are found in the reference solution but not found in the TOP ranked classes.

In table 3.3 are presented the ROC-AUC values for different attributes computed for the systems Ant, Tomcat Catalina, and Hibernate.

Table 3.3: ROC-AUC metric values extracted.

Metrics	Ant	Tomcat Catalina	Hibernate
PR_U2.W	0.95823	0.92341	0.95823
PR	0.94944	0.92670	0.94944
PR_U	0.95060	0.93220	0.95060
CONN_TOTAL.W	0.94437	0.92595	0.94437
CONN_TOTAL	0.94630	0.93903	0.94630

3.2.5 Measurements using logical dependencies

To evaluate the results obtained using logical dependencies, we used the same tool used in section X. Previously, the tool used only structural dependencies extracted from the source code of the software systems. In this chapter, we intend to add also the logical dependencies from the versioning system to observe if the results could be improved or not.

For this, the logical dependencies used were filtered based on the update percentage of the entities involved. We define a logical dependency as a connection observed via commits in the versioning system between entity A and entity B. The update percentage of entity A with entity B is determined as follows: the percentage from the total commits involving A that contains entity B.

$$\text{update percentage for } A = \frac{100 * \text{commits involving } A \text{ and } B}{\text{total nr of commits involving } A}$$

$$\text{update percentage for } B = \frac{100 * \text{commits involving } A \text{ and } B}{\text{total nr of commits involving } B}$$

We calculated the update percentage for each side of the connection (LD) and filtered the connections found based on it. The rule set is that both entities had to have an update percentage greater than the threshold value. In tables 3.4, 3.5, and 3.6, we introduced the logical dependencies among structural dependencies. We started with logical dependencies that have a percentage of update greater than 10%, which means that in at least 10% of the commits involving A or B, A and B update together. Then we increased the threshold value by 10 until we remained only with entities that update in all the commits together.

As for the new results obtained, in tables 3.4, 3.5, and 3.6, highlighted with orange, are the values that are close to the previously registered values but did not surpass them. Highlighted with green are values that are better than the previously registered values.

In tables 3.7, 3.8, and 3.9, we only used logical dependencies. The measurements obtained by using only logical dependencies are not as good as using both combined or using only structural dependencies. One possible cause is that key classes have a better design than the rest of the classes and this can make them less prone to

Table 3.4: Measurements for Ant using structural and logical dependencies combined

Metrics	$\geq 10\%$	$\geq 20\%$	$\geq 30\%$	$\geq 40\%$	$\geq 50\%$	$\geq 60\%$	$\geq 70\%$	$\geq 80\%$	$\geq 90\%$	$\geq 100\%$	Previous
PR_U2_W	0.924	0.925	0.926	0.927	0.927	0.927	0.929	0.928	0.928	0.928	0.929
PR	0.914	0.854	0.851	0.866	0.876	0.882	0.887	0.854	0.852	0.852	0.855
PR_U	0.910	0.930	0.933	0.933	0.935	0.934	0.939	0.933	0.933	0.933	0.933
CONN_T_W	0.924	0.928	0.931	0.932	0.933	0.934	0.936	0.934	0.934	0.934	0.934
CONN_T	0.840	0.886	0.904	0.909	0.915	0.923	0.932	0.935	0.936	0.936	0.942

Table 3.5: Measurements for Tomcat using structural and logical dependencies combined

Metrics	$\geq 10\%$	$\geq 20\%$	$\geq 30\%$	$\geq 40\%$	$\geq 50\%$	$\geq 60\%$	$\geq 70\%$	$\geq 80\%$	$\geq 90\%$	$\geq 100\%$	Previous
PR_U2_W	0.912	0.915	0.922	0.923	0.924	0.924	0.923	0.924	0.924	0.924	0.923
PR	0.808	0.785	0.812	0.839	0.844	0.851	0.853	0.857	0.857	0.857	0.927
PR_U	0.912	0.920	0.931	0.932	0.933	0.933	0.933	0.932	0.932	0.932	0.932
CONN_T_W	0.918	0.921	0.924	0.926	0.926	0.926	0.926	0.926	0.926	0.926	0.926
CONN_T	0.877	0.913	0.932	0.937	0.937	0.938	0.938	0.938	0.938	0.938	0.939

Table 3.6: Measurements for Hibernate using structural and logical dependencies combined

Metrics	$\geq 10\%$	$\geq 20\%$	$\geq 30\%$	$\geq 40\%$	$\geq 50\%$	$\geq 60\%$	$\geq 70\%$	$\geq 80\%$	$\geq 90\%$	$\geq 100\%$	Previous
PR_U2_W	0.955	0.957	0.958	0.958	0.958	0.958	0.958	0.958	0.958	0.958	0.958
PR	0.931	0.930	0.936	0.940	0.940	0.946	0.946	0.946	0.946	0.946	0.949
PR_U	0.942	0.946	0.948	0.949	0.949	0.950	0.950	0.950	0.950	0.950	0.951
CONN_T_W	0.939	0.942	0.943	0.944	0.944	0.944	0.945	0.945	0.945	0.945	0.944
CONN_T	0.925	0.933	0.938	0.940	0.941	0.944	0.944	0.944	0.944	0.944	0.946

change.

Table 3.7: Measurements for Ant using only logical dependencies

Metrics	$\geq 10\%$	$\geq 20\%$	$\geq 30\%$	$\geq 40\%$	$\geq 50\%$	$\geq 60\%$	$\geq 70\%$	$\geq 80\%$	$\geq 90\%$	$\geq 100\%$	Previous
PR_U2_W	0.655	0.611	0.650	0.645	0.729	0.797	0.855	0.882	0.865	0.865	0.929
PR	0.655	0.611	0.650	0.645	0.729	0.797	0.855	0.882	0.865	0.865	0.855
PR_U	0.655	0.611	0.650	0.645	0.729	0.797	0.855	0.882	0.865	0.865	0.933
CONN_T_W	0.646	0.523	0.617	0.657	0.722	0.785	0.845	0.878	0.865	0.865	0.934
CONN_T	0.646	0.523	0.617	0.657	0.722	0.785	0.845	0.878	0.865	0.865	0.942

Table 3.8: Measurements for Tomcat using only logical dependencies

Metrics	$\geq 10\%$	$\geq 20\%$	$\geq 30\%$	$\geq 40\%$	$\geq 50\%$	$\geq 60\%$	$\geq 70\%$	$\geq 80\%$	$\geq 90\%$	$\geq 100\%$	Previous
PR_U2_W	0.702	0.627	0.627	0.712	0.741	0.775	0.786	0.796	0.796	0.796	0.923
PR	0.675	0.617	0.627	0.712	0.741	0.775	0.786	0.796	0.796	0.796	0.927
PR_U	0.675	0.618	0.627	0.712	0.741	0.775	0.786	0.796	0.796	0.796	0.932
CONN_T_W	0.676	0.597	0.624	0.712	0.741	0.775	0.786	0.796	0.796	0.796	0.926
CONN_T	0.638	0.585	0.624	0.712	0.741	0.775	0.786	0.796	0.796	0.796	0.939

Table 3.9: Measurements for Hibernate using only logical dependencies

Metrics	$\geq 10\%$	$\geq 20\%$	$\geq 30\%$	$\geq 40\%$	$\geq 50\%$	$\geq 60\%$	$\geq 70\%$	$\geq 80\%$	$\geq 90\%$	$\geq 100\%$	Previous
PR_U2_W	0.651	0.596	0.597	0.616	0.619	0.644	0.649	0.650	0.650	0.650	0.958
PR	0.641	0.594	0.597	0.616	0.619	0.644	0.649	0.650	0.650	0.650	0.949
PR_U	0.641	0.595	0.597	0.616	0.619	0.644	0.649	0.650	0.650	0.650	0.951
CONN_T_W	0.652	0.591	0.597	0.616	0.619	0.644	0.649	0.650	0.650	0.650	0.944
CONN_T	0.649	0.591	0.597	0.616	0.619	0.644	0.649	0.650	0.650	0.650	0.946

3.3 Comparison of the extracted data with fan-in and fan-out metric

Fan-in and fan-out are coupling metrics. The fan-in of entity A is the total number of modules that call functions of A. The fan-out of A is the total number of entities called by A [23]. Related to fan-in and fan-out we have extracted CONN_IN and CONN_OUT. (aici mai am de explicat)

In tables 3.10, 3.11, and 3.12 we can find the metrics details for each documented key class.

In tables 3.13, 3.14, and 3.15 we can find the top 10 'best ranked' logical dependencies. As we can observe, the entities have only few structural connections.

Highlighted with orange are the key classes found in top 10. To be continued....

Table 3.10: Measurements for Ant key classes

Nr.	Classname	CONN_IN	CONN_OUT	CONN_TOTAL	LD
1	Project	191	191	214	157
2	Target	28	28	34	78
3	UnknownElement	17	17	30	90
4	RuntimeConfigurable	11	11	19	118
5	IntrospectionHelper	18	18	42	143
6	Main	1	1	14	82
7	TaskContainer	11	11	12	21
8	ProjectHelper2\$ElementHandler	1	1	13	30
9	Task	110	110	117	88
10	ProjectHelper	16	16	24	101

Table 3.11: Measurements for Tomcat Catalina key classes.

Nr.	Classname	CONN_IN	CONN_OUT	CONN_TOTAL	LD
1	Context	74	8	82	126
2	Request	48	28	76	215
3	Container	51	8	59	64
4	Response	38	12	50	90
5	StandardContext	11	38	49	216
6	Connector	23	9	32	89
7	Session	29	2	31	28
8	Valve	29	2	31	19
9	Wrapper	29	1	30	36
10	Manager	25	3	28	31
11	Host	26	1	27	44
12	Service	20	6	26	51
13	Engine	23	2	25	1
14	Realm	18	6	24	21
15	CoyoteAdapter	1	22	23	140
16	StandardHost	8	15	23	88
17	LifecycleListener	21	1	22	3
18	StandardEngine	2	19	21	57
19	Pipeline	19	2	21	20
20	Server	16	4	20	49
21	HostConfig	3	15	18	79
22	StandardWrapper	5	13	18	92
23	StandardService	3	12	15	81
24	Catalina	2	13	15	94
25	Loader	14	1	15	18
26	StandardServer	2	12	14	94
27	StandardPipeline	1	10	11	62
28	Bootstrap	3	3	6	41

Table 3.12: Measurements for Hibernate key classes.

Nr.	Classname	CONN_IN	CONN_OUT	CONN_TOTAL	LD
1	SessionFactoryImplementor	438	43	481	51
2	Type	444	5	449	0
3	Table	89	29	118	82
4	SessionImplementor	52	12	64	14
5	Criteria	45	12	57	15
6	Column	46	10	56	20
7	Session	31	21	52	52
8	Query	12	28	40	0
9	Configuration	1	38	39	115
10	SessionFactory	24	12	36	33
11	Criterion	30	3	33	0
12	Projection	11	3	14	0
13	ConnectionProvider	12	2	14	0
14	Transaction	11	1	12	0

Table 3.13: Top 10 measurements for Ant.

Nr.	Classname	CONN_IN	CONN_OUT	CONN_TOTAL	LD
1	Project	191	23	214	157
2	Project\$AntRefTable	1	2	3	157
3	Path	39	13	52	147
4	Path\$PathElement	3	2	5	147
5	IntrospectionHelper	18	24	42	143
6	IntrospectionHelper\$AttributeSetter	8	1	9	143
7	IntrospectionHelper\$Creator	3	5	8	143
8	IntrospectionHelper\$NestedCreator	7	1	8	143
9	Ant	2	15	17	136
10	Ant\$Reference	3	1	4	136

Table 3.14: Top 10 measurements for Tomcat Catalina.

Nr.	Classname	CONN_IN	CONN_OUT	CONN_TOTAL	LD
1	StandardContext	11	38	49	216
2	StandardContext\$ContextFilterMaps	0	0	0	216
3	StandardContext\$NoPluggabilityServletContext	0	0	0	216
4	Request	48	28	76	215
5	Request\$SpecialAttributeAdapter	0	0	0	215
6	ApplicationContext	3	22	25	158
7	ApplicationContext\$DispatchData	0	0	0	158
8	ContextConfig	3	26	29	143
9	ContextConfig\$DefaultWebXmlCacheEntry	0	0	0	143
10	ContextConfig\$JavaClassCacheEntry	0	0	0	143

Table 3.15: Top 10 measurements for Hibernate.

Nr.	Classname	CONN_IN	CONN_OUT	CONN_TOTAL	LD
1	AvailableSettings	1	0	1	205
2	AbstractEntityPersister	9	143	152	190
3	AbstractEntityPersister\$CacheEntryHelper	0	0	0	190
4	AbstractEntityPersister\$InclusionChecker	0	0	0	190
5	AbstractEntityPersister\$NoopCacheEntryHelper	0	0	0	190
6	AbstractEntityPersister\$ReferenceCacheEntryHelper	0	0	0	190
7	AbstractEntityPersister\$StandardCacheEntryHelper	0	0	0	190
8	AbstractEntityPersister\$StructuredCacheEntryHelper	0	0	0	190
9	Dialect	265	104	369	176
10	SessionFactoryImpl\$SessionBuilderImpl	1	25	26	167

Bibliography

- [1] Nemitari Ajenka and Andrea Capiluppi. Understanding the interplay between the logical and structural coupling of software classes. *Journal of Systems and Software*, 134:120–137, 2017.
- [2] Nemitari Ajenka, Andrea Capiluppi, and Steve Counsell. An empirical study on the interplay between semantic coupling and co-change of software classes. *Empirical Software Engineering*, 23(3):1791–1825, 2018.
- [3] G. Bavota, B. Dit, R. Oliveto, M. Di Penta, D. Poshyvanyk, and A. De Lucia. An empirical study on the developers’ perception of software coupling. In *2013 35th International Conference on Software Engineering (ICSE)*, pages 692–701, May 2013.
- [4] Fabian Beck and Stephan Diehl. On the congruence of modularity and code coupling. In *Proceedings of the 19th ACM SIGSOFT Symposium and the 13th European Conference on Foundations of Software Engineering, ESEC/FSE ’11*, pages 354–364, New York, NY, USA, 2011. ACM.
- [5] David Binkley. Source code analysis: A road map. pages 104–119, 06 2007.

- [6] Grady Booch. *Object-Oriented Analysis and Design with Applications (3rd Edition)*. Addison Wesley Longman Publishing Co., Inc., Redwood City, CA, USA, 2004.
- [7] Andrew P. Bradley. The use of the area under the roc curve in the evaluation of machine learning algorithms. *Pattern Recognition*, 30(7):1145–1159, 1997.
- [8] Trosky B. Callo Arias, Pieter van der Spek, and Paris Avgeriou. A practice-driven systematic review of dependency analysis solutions. *Empirical Software Engineering*, 16(5):544–586, Oct 2011.
- [9] M. Cataldo, A. Mockus, J. A. Roberts, and J. D. Herbsleb. Software dependencies, work dependencies, and their impact on failures. *IEEE Transactions on Software Engineering*, 35(6):864–878, Nov 2009.
- [10] Marcelo Cataldo, Audris Mockus, Jeffrey A. Roberts, and James D. Herbsleb. Software dependencies, work dependencies, and their impact on failures. *IEEE Transactions on Software Engineering*, 35:864–878, 2009.
- [11] M. L. Collard, H. H. Kagdi, and J. I. Maletic. An XML-based lightweight C++ fact extractor. In *Proceedings of the 11th IEEE International Workshop on Program Comprehension, IWPC '03*, pages 134–, Washington, DC, USA, 2003. IEEE Computer Society.
- [12] Michael L. Collard, Michael J. Decker, and Jonathan I. Maletic. Lightweight transformation and fact extraction with the srcML toolkit. In *Proceedings of the 2011 IEEE 11th International Working Conference on Source Code Analysis and Manipulation, SCAM '11*, pages 173–184, Washington, DC, USA, 2011. IEEE Computer Society.

- [13] Ben Collins-Sussman, Brian W. Fitzpatrick, and C. Michael Pilato. *Version Control With Subversion for Subversion 1.6: The Official Guide And Reference Manual*. CreateSpace, Paramount, CA, 2010.
- [14] Ioana Şora. Helping program comprehension of large software systems by identifying their most important classes. In *Evaluation of Novel Approaches to Software Engineering - 10th International Conference, ENASE 2015, Barcelona, Spain, April 29-30, 2015, Revised Selected Papers*, pages 122–140. Springer International Publishing, 2015.
- [15] Yi Ding, B. Li, and Peng He. An improved approach to identifying key classes in weighted software network. *Mathematical Problems in Engineering*, 2016:1–9, 2016.
- [16] L. do Nascimento Vale and M. de A. Maia. Keele: Mining key architecturally relevant classes using dynamic analysis. In *2015 IEEE International Conference on Software Maintenance and Evolution (ICSME)*, pages 566–570, 2015.
- [17] Tom Fawcett. An introduction to roc analysis. *Pattern Recognition Letters*, 27(8):861–874, 2006. ROC Analysis in Pattern Recognition.
- [18] Harald Gall, Karin Hajek, and Mehdi Jazayeri. Detection of logical coupling based on product release history. In *Proceedings of the International Conference on Software Maintenance, ICSM '98*, pages 190–, Washington, DC, USA, 1998. IEEE Computer Society.
- [19] Harald Gall, Mehdi Jazayeri, and Jacek Krajewski. Cvs release history data for detecting logical couplings. In *Proceedings of the 6th International Workshop*

on *Principles of Software Evolution*, IWPSE '03, pages 13–, Washington, DC, USA, 2003. IEEE Computer Society.

- [20] Eirini Kalliamvakou, Georgios Gousios, Kelly Blincoe, Leif Singer, Daniel M. German, and Daniela Damian. An in-depth study of the promises and perils of mining github. *Empirical Software Engineering*, 21(5):2035–2071, Oct 2016.
- [21] M. Kamran, M. Ali, and B. Akbar. Identification of core architecture classes for object-oriented software systems. *Journal of Applied Computer Science & Mathematics*, 10:21–25, 2016.
- [22] P. Meyer, H. Siy, and S. Bhowmick. Identifying important classes of large software systems through k-core decomposition. *Adv. Complex Syst.*, 17, 2014.
- [23] A. Mubarak, S. Counsell, and R. M. Hierons. An evolutionary study of fan-in and fan-out metrics in oss. In *2010 Fourth International Conference on Research Challenges in Information Science (RCIS)*, pages 473–482, 2010.
- [24] Gustavo Ansal di Oliva and Marco Aurelio Gerosa. On the interplay between structural and logical dependencies in open-source software. In *Proceedings of the 2011 25th Brazilian Symposium on Software Engineering, SBES '11*, pages 144–153, Washington, DC, USA, 2011. IEEE Computer Society.
- [25] Ioana Șora. Finding the right needles in hay - helping program comprehension of large software systems. In *Proceedings of the 10th International Conference on Evaluation of Novel Approaches to Software Engineering - Volume 1: ENASE,*, pages 129–140. INSTICC, SciTePress, 2015.

- [26] Ioana Șora and Ciprian-Bogdan Chirila. Finding key classes in object-oriented software systems by techniques based on static analysis. *Information and Software Technology*, 116:106176, 2019.
- [27] M. H. Osman, M. R. V. Chaudron, and P. v. d. Putten. An analysis of machine learning algorithms for condensing reverse engineered class diagrams. In *2013 IEEE International Conference on Software Maintenance*, pages 140–149, 2013.
- [28] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The pagerank citation ranking: Bringing order to the web. Technical Report 1999-66, Stanford InfoLab, November 1999. Previous number = SIDL-WP-1999-0120.
- [29] Weifeng Pan, Beibei Song, Kangshun Li, and Kejun Zhang. Identifying key classes in object-oriented software using generalized k-core decomposition. *Future Generation Computer Systems*, 81:188–202, 2018.
- [30] Neeraj Sangal, Ev Jordan, Vineet Sinha, and Daniel Jackson. Using dependency models to manage complex software architecture. In *Proceedings of the 20th Annual ACM SIGPLAN Conference on Object-oriented Programming, Systems, Languages, and Applications*, OOPSLA '05, pages 167–176, New York, NY, USA, 2005. ACM.
- [31] Ioana Șora. A PageRank based recommender system for identifying key classes in software systems. In *2015 IEEE 10th Jubilee International Symposium on Applied Computational Intelligence and Informatics (SACI)*, pages 495–500, May 2015.
- [32] Adelina Diana Stana. and Ioana Șora. Identifying logical dependencies from co-changing classes. In *Proceedings of the 14th International Conference on*

Evaluation of Novel Approaches to Software Engineering - Volume 1: ENASE, pages 486–493. INSTICC, SciTePress, 2019.

- [33] Stana Adelina and Șora Ioana. Analyzing information from versioning systems to detect logical dependencies in software systems. In *International Symposium on Applied Computational Intelligence and Informatics (SACI)*, May 2019.
- [34] D. Steidl, B. Hummel, and E. Juergens. Using network analysis for recommendation of central software classes. In *2012 19th Working Conference on Reverse Engineering*, pages 93–102, 2012.
- [35] L. Tahvildari and K. Kontogiannis. Improving design quality using meta-pattern transformations: a metric-based approach. *J. Softw. Maintenance Res. Pract.*, 16:331–361, 2004.
- [36] Ferdian Thung, David Lo, Mohd Hafeez Osman, and Michel R. V. Chaudron. Condensing class diagrams by analyzing design and network metrics using optimistic classification. In *Proceedings of the 22nd International Conference on Program Comprehension, ICPC 2014*, page 110–121, New York, NY, USA, 2014. Association for Computing Machinery.
- [37] X. Yang, D. Lo, X. Xia, and J. Sun. Condensing class diagrams with minimal manual labeling cost. In *2016 IEEE 40th Annual Computer Software and Applications Conference (COMPSAC)*, volume 1, pages 22–31, 2016.
- [38] Liguu Yu. Understanding component co-evolution with a study on linux. *Empirical Softw. Engg.*, 12(2):123–141, April 2007.
- [39] A. Zaidman, T. Calders, S. Demeyer, and J. Paredaens. Applying webmining techniques to execution traces to support the program comprehension process. In

Ninth European Conference on Software Maintenance and Reengineering, pages 134–142, 2005.

- [40] Andy Zaidman and Serge Demeyer. Automatic identification of key classes in a software system using webmining techniques. *Journal of Software Maintenance and Evolution: Research and Practice*, 20(6):387–417, 2008.