

Logical dependencies: extraction from the versioning system and an example of usage *

Adelina Stana¹ and Ioana Sora²

¹ Stana Adelina Diana

Address

stana.adelina.diana@gmail.com

² Sora Ioana

Address

ioana.sora@cs.upt.ro

Abstract. One of the biggest problems when it comes to legacy software systems is the lack of modernization. A legacy software system is an obsolete system that might still be in use and satisfy the organization's needs but due to the lack of modernization, it has a high maintenance cost. One solution is to modernize or refactor the system and to do that the developers must have a good understanding of the system. Here comes one other problem of the legacy software systems, the lack of up-to-date documentation or the lack of documentation. The solution for this problem is to analyze the legacy software system in order to gain more knowledge about the system and how the system works. Our approach is to analyze the legacy software system by using historical information extracted from the versioning systems.

Keywords: logical dependencies, versioning system, key classes.



A *logical dependency* between two software entities (classes, modules, interfaces, etc.) signals that those entities in some way depend on one another.

The concept of logical coupling (dependency) was first introduced by Gall [7]. They defined the logical coupling between two modules as the fact that the modules repeatedly change together during the historical evolution of the software system.

To define the logical dependencies, we have first to define the co-changing pairs. A *co-changing pair* are two software entities that update together in the same commit. For example, a commit that contains seven entities will generate 21 co-changing pairs ($C_k^n = \frac{n!}{k!(n-k)!} = \frac{7!}{2!(5)!} = 21$).

For a co-changing pair to be considered a logical dependency, it has to pass multiple filters. The filter's purpose is to increase the confidence that the remaining co-changing pairs are indeed related and decrease the data volume. In our previous work, we applied multiple filters to the co-changing pairs extracted and observed how the filtering affects the number of co-changing pairs. From the filtering methods previously explored, we choose to keep the commit size filter [16]. In addition, we added a new filter called the connection strength filter.

The commit size filter filters out all co-changing pairs relating to commits with more than 10 files changed. Our previous works state that commits with more than 10 files

* If this is an extended version of a conference paper, it should be clearly stated here.

1 changed tend to be code unrelated; in addition, only a tiny percent of the total commits
 2 have more than 10 files changed; so, the filtering will take out only a few commits.

3 After applying the commit size filtering, we apply the filtering based on connection
 4 strength. To determine the connection strength, we first need to calculate the connection
 5 factors for both entities that form a co-changing pair.

6 Assuming that we have a co-changing pair formed by entities A and B, the connection
 7 factor of entity A with entity B is the percentage from the total commits involving A that
 8 contains entity B. The connection factor of entity B with entity A is the percentage from
 9 the total commits involving B that contain also entity A.

$$\text{connection factor for A} = \frac{100 * \text{commits involving A and B}}{\text{total nr of commits involving A}} \quad (1)$$

$$\text{connection factor for B} = \frac{100 * \text{commits involving A and B}}{\text{total nr of commits involving B}} \quad (2)$$

10 We calculated the connection factor for each entity involved in a co-changing pair and
 11 filtered the co-changing pairs based on it. The rule set is that both entities had to have a
 12 connection factor with each other greater than the threshold value.

13 After the filtering the remaining co-changing pairs, now called logical dependen-
 14 cies, are exported in CSV files.

15 The entire process of extracting co-changing pairs from the versioning system, filter
 16 them, and export the remaining ones into CSV files is done with a tool written in Python.

17 The workflow is presented in figure 1

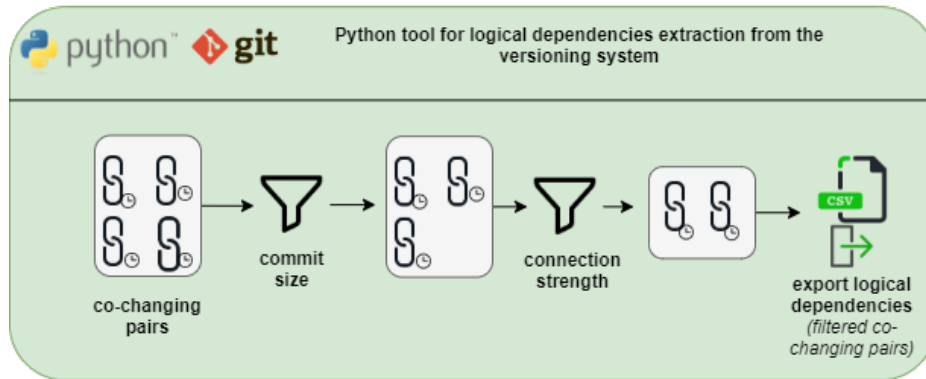


Fig. 1. Workflow for logical dependencies detection

18 1. The concept of key classes

19 Zaidman et al [22] were the first to introduce the concept of key classes and it refers to
 20 classes that can be found in documents written to provide an architectural overview of
 21 the system or an introduction to the system structure. Tahvildari and Kontogiannis have a

1 more detailed definition regarding key classes concept: “Usually, the most important con-
 2 cepts of a system are implemented by very few key classes which can be characterized
 3 by the specific properties. These classes, which we refer to as key classes, manage many
 4 other classes or use them in order to implement their functionality. The key classes are
 5 tightly coupled with other parts of the system. Additionally, they tend to be rather com-
 6 plex, since they implement much of the legacy system’s functionality” [18]. Also, other
 7 researchers use a similar concept as the one defined by Zaidman but under different terms
 8 like important classes [9] or central software classes [17].

9 The key class identification can be done by using different algorithms with different
 10 inputs. In the research of Osman et al., the key class identification is made by using a
 11 machine learning algorithm and class diagrams as input for the algorithm [12]. Thung
 12 et al. builds on top of Osman et al.’s approach and adds network metrics and optimistic
 13 classification in order to detect key classes [19].

14 Zaidman et al. use a webmining algorithm and dynamic analysis of the source code to
 15 identify the key classes [22].

16 Sora et al. use a page ranking algorithm for finding key classes and static analysis of
 17 the source code [2], [10], [3], [15]. In [11] the authors use in addition to the previous re-
 18 search also other class attributes to identify important classes. The page ranking algorithm
 19 is a customization of PageRank, the algorithm used to rank web pages [13]. The PageR-
 20 ank algorithm works based on a recommendation system. If one node has a connection
 21 with another node, then it recommends the second node. In previous works, connections
 22 are established based on structural dependencies extracted from static code analysis. If A
 23 has a structural dependency with B, then A recommends B, and also B recommends A.

24 The ranking algorithm ranks all the classes from the source code of the system ana-
 25 lyzed according to their importance. To identify the important classes from the rest of the
 26 classes a threshold for TOP classes from the top of the ranking is set. The TOP threshold
 27 value can go from 1 to the total number of classes found in the system.

28 Some researchers [22], [4], [14] consider that 15% of the total number of classes of
 29 the system is a suited value for the TOP threshold. Other researchers [11] consider that
 30 15% of the total number of classes is a too high value for the TOP threshold and suggest
 31 that a value in the range of 20–30 is better.

32 2. Metrics for results evaluation

33 To evaluate the quality of the key classes ranking algorithm and solution produced, the
 34 key classes found by the algorithm are compared with a reference solution.

35 The reference solution is extracted from the developer documentation. Classes men-
 36 tioned in the documentation are considered key classes and form the reference solution
 37 (ground truth) used for validation [20].

38 For the comparison between both solutions, a classification model is used. The quality
 39 of the solution produced is evaluated by using metrics that evaluate the performance of
 40 the classification model, such as Precision-Recall and Receiver Operating Characteristic
 41 Area Under Curve (ROC-AUC).

42 A classification model (or “classifier”) is a mapping between expected results and
 43 predicted results [6], [1].

44 Precision is the ratio of True Positives to all the positives of the result set.

$$precision = \frac{TP}{TP + FN} \quad (3)$$

1 The recall is the ratio of True Positives to all the positives of the reference set.

$$recall = \frac{TP}{TP + FP} \quad (4)$$

2 As mentioned in section 1, to distinguish the key classes from the rest of the classes a
3 TOP threshold is used. Some researchers consider that 15% of the total classes is the best
4 value for the TOP threshold and others consider that the value should be in the range of
5 20-30.

6 The precision-recall metric is suited if the threshold value is fixed. If the threshold
7 value is variable, then metrics that capture the behavior over all possible values must be
8 used. Such metric is the Receiver Operating Characteristic metric.

9 *Receiver Operating Characteristic Area Under Curve*

10 The ROC graph is a two-dimensional graph that has on the X-axis plotted the false
11 positive rate and on the Y-axis the true positive rate. By plotting the true positive rate and
12 the false positive rate at thresholds that vary between a minimum and a maximum possible
13 value we obtain the ROC curve. The area under the ROC curve is called Area Under the
14 Curve (AUC).

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

$$True\ positive\ rate(TPR) = \frac{TP}{TP + FN} \quad (5)$$

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

$$False\ positive\ rate(FPR) = \frac{FP}{FP + TN} \quad (6)$$

19 In multiple related works, the ROC-AUC metric has been used to evaluate the results
20 for finding key classes of software systems. For a classifier to be considered good, its
21 ROC-AUC metric value should be as close to 1 as possible, when the value is 1 then the
22 classifier is considered to be perfect.

23 Osman et al. obtained in their research an average Area Under the Receiver Operating
24 Characteristic Curve (ROC-AUC) score of 0.750 [12]. Thung et al. obtained an average
25 ROC-AUC score of 0.825 [19] and Sora et al. obtained an average ROC-AUC score of
26 0.894 [11].

27 3. Baseline versus current approach



28 We use the research of I. Sora et al [11] as a baseline for our research involving the usage
29 of logical dependencies to find key classes. The baseline approach uses a tool that takes
30 as an input the source code of the system and applies ranking strategies to rank the classes
31 according to their importance.

In order to rank the classes according to their importance, different class metrics are used [4], [22], [14]. Below are some of the class metrics used in the baseline approach and our current research to rank the classes according to their importance. We use only a subset of the metrics used in previous research because the extracted logical dependencies are undirected. The class metrics used can be separated into two categories: class connection metrics and class Pagerank values.

The class connection metrics are CONN-TOTAL-W, which is the total weight of all connections of the class, and CONN-TOTAL, the total number of distinct classes that a class uses or are used by a class [11].

Previous research used PageRank values computed on both directed and undirected, weighted and unweighted graphs. In the current research, we used the PR, which is the PageRank value computed on the directed and unweighted graph. The PR-U, which is the value computed on the undirected and unweighted graph, and PR-U2-W, the value computed on the weighted graph with back-recommendations [2], [10], [11], [15].

Based on the class attributes presented, all the classes of the system are ranked. To differentiate the important (key) classes from the rest of the classes, a TOP threshold for the top classes found is set. The threshold vary between 20 and 30 classes.

The baseline approach not only identifies the key classes but also evaluates the performance of the solution produced. The same approach as the one presented in section 2 is used for the evaluation of the results. The key classes found by the ranking algorithm are compared with a reference solution that is extracted from the developer documentation by using a classification model.

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.

Due to the fact that the TOP threshold is varied, the Receiver Operating Characteristic Area Under Curve metric is used for the evaluation of the results.

We modified the tool to take the logical dependencies as input; the rest of the workflow remained the same as the baseline approach. The entire workflow is presented in figure 2.

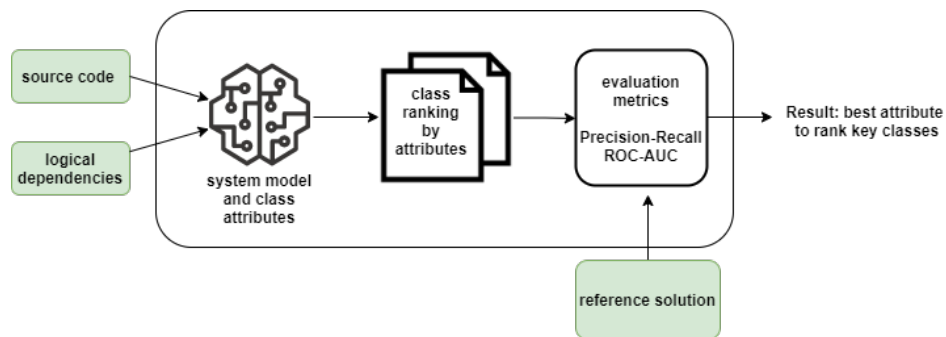


Fig. 2. Overview of the current approach.

4. Data set used

In this section, we will look over all the systems studied in the baseline research presented in section 3, and we will try to identify the systems that could be used also in our current research involving logical dependencies.

The research of I. Sora et al [11] takes into consideration structural public dependencies that are extracted using static analysis techniques and was performed on the object-oriented systems presented in table 1.

The requirements for a system to qualify as suited for investigations using logical dependencies are: has to be on GitHub, has to have releases to identify a specific version (previous research was done only on specific releases), and also, has to have a significant number of commits. From the total of 14 object-oriented systems listed in the paper [11], 13 of them have repositories in Github 1. And from the found repositories, only 6 repositories have the same release tag as the specified version in previous research. 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 1 Apache Ant is the most used and analyzed in other works [16], [5], [21], [8].

Table 1. 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

5. Measurements using logical dependencies

As we mentioned in the beginning the purpose is to check if the logical dependencies can improve key class detection.

As presented in section 3, and section 1 the key class detection was done by using structural dependencies of the system. In this section, we will use the same tool used in the baseline approach presented in section 3, and we will add a new input to it, the logical dependencies.

Below is a comparison between the new approach and baseline approach, how we collect the logical dependencies, the results obtained previously, and the new results obtained. The new results are separated into two categories, the results obtained by using structural and logical dependencies and the results obtained by using only logical dependencies.

5.1. Measurements using only the baseline approach

In table 2 are presented the ROC-AUC values for different attributes computed for the systems Ant, Tomcat Catalina, and Hibernate by using the baseline approach. We intend to compare these values with the new values obtained by using also logical dependencies in key class detection.

Table 2. 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

5.2. Measurements using combined structural and logical dependencies

The tool used in the baseline approach runs a graph-ranking algorithm. The graph used contains the structural dependencies extracted from static source code analysis. Each edge in the graph represents a dependency, the entities that form a structural dependency are represented as vertices in the graph. As mentioned in section ??, we modified the tool to read also logical dependencies and add them to the graph. In this section, we add in the graph the logical dependencies together with the structural dependencies.

In tables 3, 4, and 5, on each line, we have the metric that is calculated and on each column, we have the connection strength threshold that was applied to the logical dependencies used in identifying the key classes. We started with logical dependencies that have a connection strength 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. The last column contains the results obtained previously by the tool by only using structural dependencies.

As for the new results obtained by combining structural and logical dependencies, 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

- 1 registered values. At this step, we can also observe that for all three systems measured in
 2 tables 3, 4, and 5, the best values obtained are for connection strength between 40-70%.

Table 3. Measurements for Ant using structural and logical dependencies combined

Metrics	$\geq 10\%$	$\geq 20\%$	$\geq 30\%$	$\geq 40\%$	$\geq 50\%$	$\geq 60\%$	$\geq 70\%$	$\geq 80\%$	$\geq 90\%$	$\geq 100\%$	Baseline
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
CON_T_W	0.924	0.928	0.931	0.932	0.933	0.934	0.936	0.934	0.934	0.934	0.934
CON_T	0.840	0.886	0.904	0.909	0.915	0.923	0.932	0.935	0.936	0.936	0.942

Table 4. Measurements for Tomcat using structural and logical dependencies combined

Metrics	$\geq 10\%$	$\geq 20\%$	$\geq 30\%$	$\geq 40\%$	$\geq 50\%$	$\geq 60\%$	$\geq 70\%$	$\geq 80\%$	$\geq 90\%$	$\geq 100\%$	Baseline
PR_U2_W	0.910	0.917	0.923	0.924	0.924	0.924	0.924	0.924	0.924	0.924	0.923
PR	0.811	0.800	0.815	0.834	0.847	0.852	0.853	0.858	0.858	0.858	0.927
PR_U	0.910	0.921	0.931	0.933	0.933	0.932	0.933	0.932	0.932	0.932	0.932
CON_T_W	0.914	0.920	0.924	0.926	0.926	0.926	0.926	0.926	0.926	0.926	0.926
CON_T	0.868	0.906	0.930	0.936	0.937	0.938	0.938	0.938	0.938	0.938	0.939

Table 5. Measurements for Hibernate using structural and logical dependencies combined

Metrics	$\geq 10\%$	$\geq 20\%$	$\geq 30\%$	$\geq 40\%$	$\geq 50\%$	$\geq 60\%$	$\geq 70\%$	$\geq 80\%$	$\geq 90\%$	$\geq 100\%$	Baseline
PR_U2_W	0.954	0.957	0.958	0.958	0.958	0.958	0.958	0.958	0.958	0.958	0.958
PR	0.929	0.929	0.933	0.939	0.939	0.946	0.947	0.947	0.947	0.947	0.949
PR_U	0.942	0.947	0.948	0.949	0.949	0.950	0.950	0.950	0.950	0.950	0.951
CON_T_W	0.939	0.942	0.943	0.944	0.944	0.945	0.945	0.945	0.945	0.945	0.944
CON_T	0.924	0.933	0.938	0.941	0.941	0.944	0.945	0.945	0.945	0.945	0.946

3 5.3. Measurements using only logical dependencies

4 In the previous section, we added in the graph based on which the ranking algorithm
 5 works the logical and structural dependencies. In the current section, we will add only the
 6 logical dependencies to the graph.

7 In tables 6, 7, and 8, are presented the results obtained by using only logical depen-
 8 dencies to detect key classes. The measurements obtained are not as good as using logical
 9 and structural dependencies combined or using only structural dependencies. But, all the
 10 values obtained are above 0.5, which means that a good part of the key classes is detected
 11 by only using logical dependencies. As mentioned in section 2, a classifier is good if it
 12 has the ROC-AUC value as close to 1 as possible.

- 1 One possible explanation for the less performing results is that the key classes may
 2 have a better design than the rest of the classes, which means that are less prone to change.
 3 If the key classes are less prone to change, this implies that the number of dependencies
 4 extracted from the versioning system can be less than for other classes.

Table 6. Measurements for Ant using only logical dependencies

Metrics	$\geq 10\%$	$\geq 20\%$	$\geq 30\%$	$\geq 40\%$	$\geq 50\%$	$\geq 60\%$	$\geq 70\%$	$\geq 80\%$	$\geq 90\%$	$\geq 100\%$	Baseline
PR_U2_W	0.720	0.627	0.718	0.703	0.732	0.824	0.852	0.881	0.876	0.876	0.929
PR	0.720	0.627	0.718	0.703	0.732	0.824	0.852	0.881	0.876	0.876	0.855
PR_U	0.720	0.627	0.718	0.703	0.732	0.824	0.852	0.881	0.876	0.876	0.933
CON_T_W	0.722	0.581	0.644	0.676	0.727	0.819	0.842	0.874	0.876	0.876	0.934
CON_T	0.722	0.581	0.644	0.676	0.727	0.819	0.842	0.874	0.876	0.876	0.942

Table 7. 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	0.672	0.656	0.645	0.697	0.754	0.776	0.786	0.799	0.799	0.799	0.923
PR	0.685	0.643	0.642	0.697	0.754	0.776	0.786	0.799	0.799	0.799	0.927
PR_U	0.685	0.643	0.644	0.697	0.754	0.776	0.786	0.799	0.799	0.799	0.932
CON_T_W	0.694	0.636	0.636	0.697	0.754	0.776	0.786	0.799	0.799	0.799	0.926
CON_T	0.654	0.611	0.636	0.697	0.754	0.776	0.786	0.799	0.799	0.799	0.939

Table 8. Measurements for Hibernate using only logical dependencies

Metrics	$\geq 10\%$	$\geq 20\%$	$\geq 30\%$	$\geq 40\%$	$\geq 50\%$	$\geq 60\%$	$\geq 70\%$	$\geq 80\%$	$\geq 90\%$	$\geq 100\%$	Baseline
PR_U2_W	0.657	0.564	0.601	0.619	0.622	0.650	0.653	0.654	0.654	0.654	0.958
PR	0.644	0.564	0.601	0.619	0.622	0.650	0.653	0.654	0.654	0.654	0.949
PR_U	0.644	0.564	0.601	0.619	0.622	0.650	0.653	0.654	0.654	0.654	0.951
CON_T_W	0.649	0.564	0.601	0.619	0.622	0.650	0.653	0.654	0.654	0.654	0.944
CON_T	0.644	0.564	0.601	0.619	0.622	0.650	0.653	0.654	0.654	0.654	0.946

5 6. Correlation between details of the systems and results

6 In this section, we discuss about the correlation between the details of the systems and the
 7 results obtained in section 5.

8 The reason why we are doing this correlation is to find if there are some links between
 9 the details of the systems and the results obtained.

10 The results obtained are presented in figures 3 - 8. We are using plots to display the
 11 results obtained to have a clearer view of how the results fluctuate over different thresholds
 12 values.

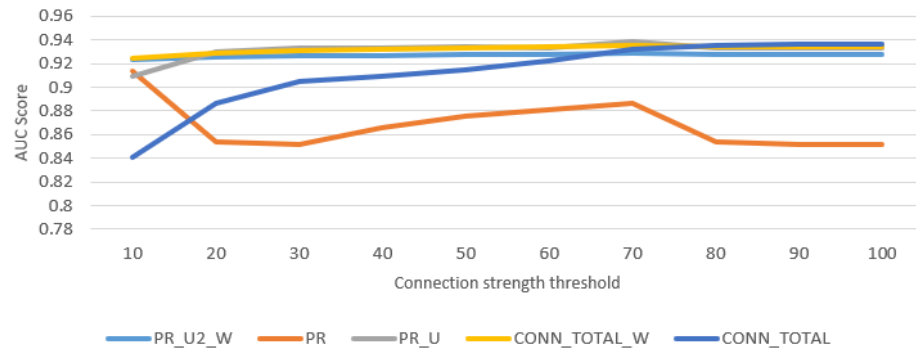


Fig. 3. Variation of AUC score when varying connection strength threshold for Ant. Results for structural and logical dependencies combined.

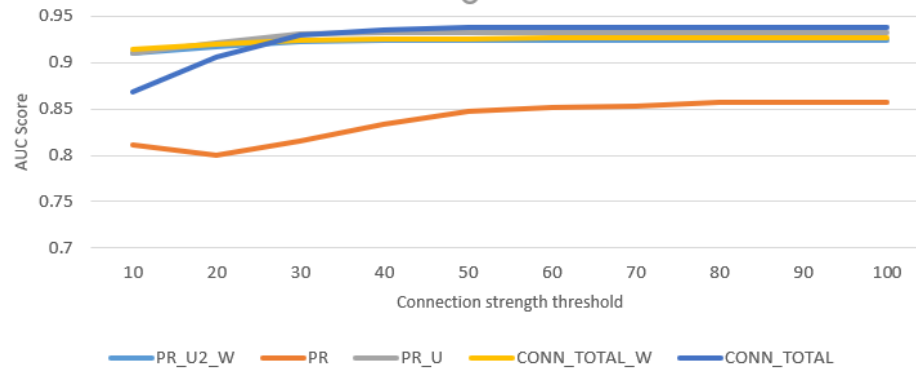


Fig. 4. Variation of AUC score when varying connection strength threshold for Tomcat. Results for structural and logical dependencies combined.

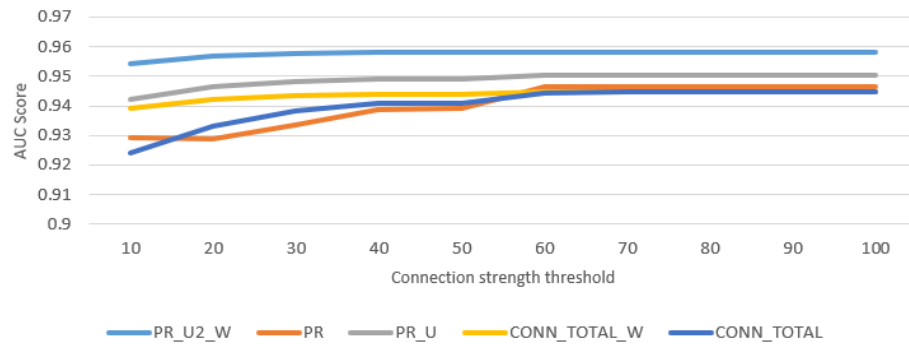


Fig. 5. Variation of AUC score when varying connection strength threshold for Hibernate. Results for structural and logical dependencies combined.

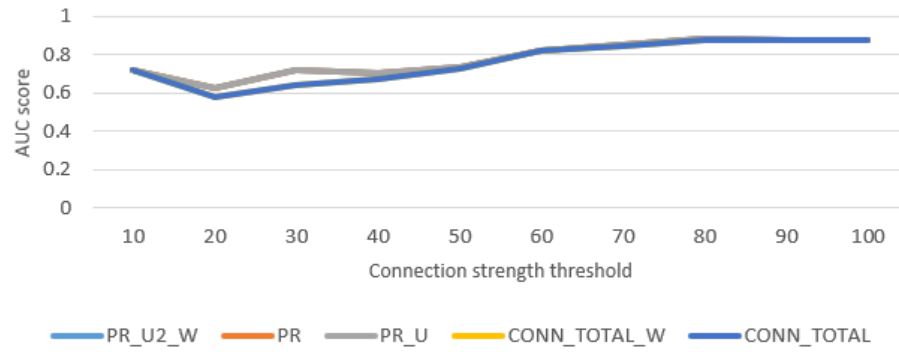


Fig. 6. Variation of AUC score when varying connection strength threshold for Ant. Results for logical dependencies only.

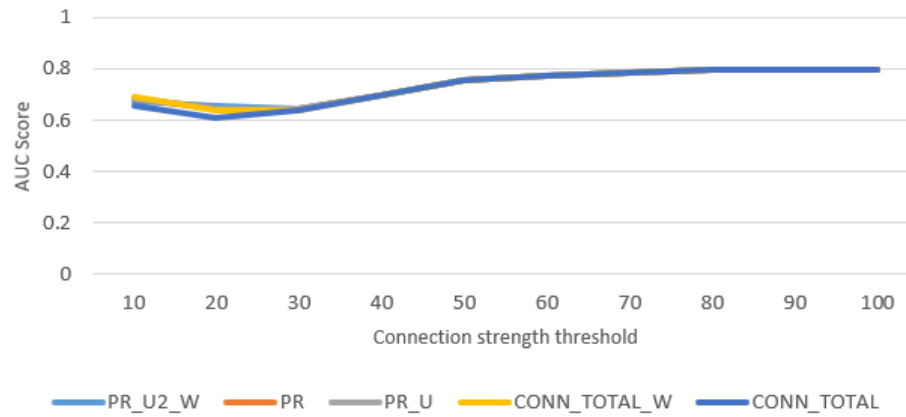


Fig. 7. Variation of AUC score when varying connection strength threshold for Tomcat. Results for logical dependencies only.

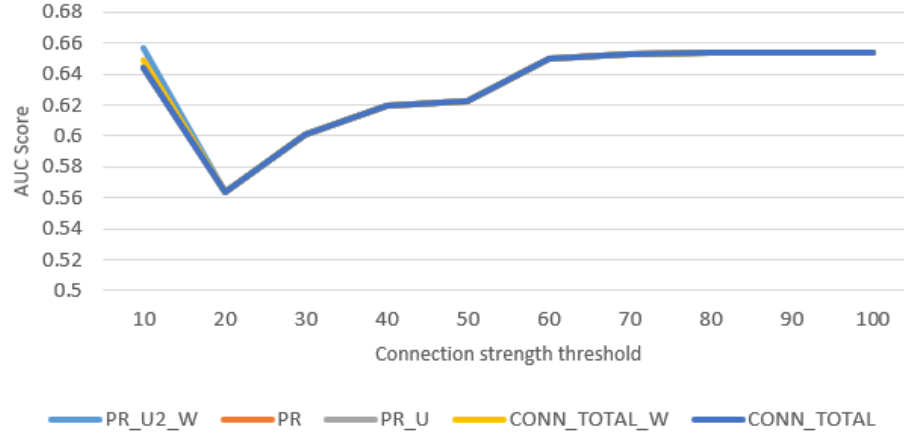


Fig. 8. Variation of AUC score when varying connection strength threshold for Hibernate. Results for logical dependencies only.

The details of the systems are presented in two tables. In table 9 are the overlappings between structural and logical dependencies expressed in percentages. Each column represents the percentage of logical dependencies that are also structural, for each column the logical dependencies are obtained by applying a different connection strength filter. The connection strength filter begins at 10, meaning that in at least 10 % of the total commits involving two entities, the entities update together. We increase the connection strength filter by 10 up until we reach 100, meaning that in all the commits that involve one entity, the other entity is present also.

In table 10 are the ratio numbers between structural dependencies and logical dependencies. We added this table in order to highlight how different the total number of both dependencies is.

Table 9. Percentage of logical dependencies that are also structural dependencies

System	$\geq 10\%$	$\geq 20\%$	$\geq 30\%$	$\geq 40\%$	$\geq 50\%$	$\geq 60\%$	$\geq 70\%$	$\geq 80\%$	$\geq 90\%$	$\geq 100\%$
Ant	25.202	34.419	36.385	34.656	33.528	33.333	28.659	33.333	35.294	35.294
Tomcat Catalina	4.059	22.089	25.000	25.758	25.926	37.525	47.368	55.285	75.000	76.923
Hibernate	6.546	26.607	29.565	32.374	32.543	45.170	44.980	42.473	42.473	42.473

Table 10. Ratio between structural and logical dependencies (SD/LD)

System	$\geq 10\%$	$\geq 20\%$	$\geq 30\%$	$\geq 40\%$	$\geq 50\%$	$\geq 60\%$	$\geq 70\%$	$\geq 80\%$	$\geq 90\%$	$\geq 100\%$
Ant	1.315	3.284	4.972	5.603	6.175	10.697	12.915	27.154	41.529	41.529
Tomcat Catalina	0.120	0.923	1.313	1.531	1.619	3.177	7.092	13.146	67.375	124.385
Hibernate	1.037	6.391	10.037	14.947	18.940	54.248	83.442	111.704	111.704	111.704

In figures 3, 4 and 5 are the measurements obtained by using structural and logical dependencies combined. In all three figures, the measurements at the beginning are smaller than the rest. Once with the increasing of the threshold value also the measurements begin to increase. Meaning that better results for key class detection are found. The best measurements are when the threshold value is between 40 and 60, after that, the measurements tend to decrease a little bit and stay at that fixed value.

A possible explanation of the results fluctuation and then capping is that if we are looking at table 10 we can see that at the beginning, the total number of logical dependencies used is close to the number of existing structural dependencies. The high volume of logical dependencies introduced might cause an erroneous detection of the key classes, in consequence, smaller measurements. When the threshold begins to be more restrictive and the total number of logical dependencies used begins to decrease, the key classes detection starts to improve. This improvement stops after the threshold value reaches 60%. If we look again at table 10 we can see that after 60% the number of structural dependencies outnumbers the number of logical dependencies up to 124 times in some cases. In addition, if we look at table 9 we can see that the remaining logical dependencies overlap a lot with the structural dependencies, so we are not introducing too much new information.

So, the number of logical dependencies used is so small that it doesn't influence the key class identification. Since the structural dependencies used don't change, we obtain the same results for different threshold values.

In figures 6, 7 and 8 are the measurements obtained by using only logical dependencies. Initially, we expected to see a Gaussian curve, but instead, we see a bell curve. We think that in the beginning, we use a high number of logical dependencies in key class detection, among those logical dependencies is an important number of key classes and also an important number of other classes. But the number of other classes does not influence the key classes detection. When we start to increase the value of the threshold and filter more the logical dependencies, we also filter some of the initial detected key classes and remain with a significant number of other classes. In this case, the other classes that remain influence the measurements, causing the worst-performing solutions. Some of the key classes are strongly connected in the versioning system, and even for higher threshold values don't get filtered out. Meanwhile, the rest of the classes that are not key classes get filtered out for higher threshold values which leads to better performing measurements when the threshold value are above 60%.

7. Conclusions

The logical dependencies are filtered co-changing pairs extracted from the versioning system history. The filters applied to the co-changing pairs are the following: the filter based on commit size and the filter based on connection strength.

In our experiments the filter based on commit size had a hard threshold of 10 files, meaning that we use co-changing pairs only from commits with at most 10 files changed. The filter based on connection strength had a variable threshold, starting with 10% and ending with 100%. We used a variable threshold for connection strength because we wanted to observe how this threshold will impact the key classes detection.

In section 5 we approached two scenarios to detect key classes by using logical dependencies. In the first scenario, we used logical dependencies together with structural

dependencies and in the second, we used only logical dependencies to detect the key classes. We modified the tool used in the baseline approach to use also logical dependencies, and then we performed the key class identification using that tool. The quality of the results obtained was evaluated with the same tool, the metric used to evaluate the results is Area Under the Receiver Operating Characteristic Curve (ROC-AUC). We then compared the evaluation results with the results obtained by the baseline approach.

Based on the results obtained, compared with the baseline results, we did saw a slight improvement in key class detection when both logical and structural dependencies were used together, the best results were obtained with a connection strength threshold of 40-70%. When we used only logical dependencies to detect key classes, the results were less performing than using only structural or structural and logical dependencies combined.

As we mentioned in section 1, also other researchers tried to identify the key classes, and even though the approaches are not the same, most of them have used the ROC-AUC metric to evaluate the quality of the results. Osman et al. obtained in their research an average Area Under the Receiver Operating Characteristic Curve (ROC-AUC) score of 0.750 [12]. Thung et al. obtained an average ROC-AUC score of 0.825 [19] and Sora et al. (the baseline approach) obtained an average ROC-AUC score of 0.894 [11].

In the current research, we obtained an average ROC-AUC score of 0.926 when using logical and structural dependencies combined and a score of 0.747 when using only logical dependencies to detect key classes.

In conclusion, by using both dependencies combined, we can obtain a slightly better ROC-AUC score than the one obtained by the baseline approach. And, by using only logical dependencies we don't obtain a better score than the baseline approach but compared with the results obtained by other researchers [12], the score obtained is almost equal. The advantage of using only logical dependencies in key class detection is that it only uses data extracted from the versioning system and can be generalized to various programming languages.

References

1. Bradley, A.P.: The use of the area under the roc curve in the evaluation of machine learning algorithms. *Pattern Recognition* 30(7), 1145–1159 (1997), <https://www.sciencedirect.com/science/article/pii/S0031320396001422>
2. Şora, I.: 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*. pp. 122–140. Springer International Publishing (2015)
3. Şora, I.: Helping program comprehension of large software systems by identifying their most important classes. In: Maciaszek, L.A., Filipe, J. (eds.) *Evaluation of Novel Approaches to Software Engineering*. pp. 122–140. Springer International Publishing, Cham (2016)
4. Ding, Y., Li, B., He, P.: An improved approach to identifying key classes in weighted software network. *Mathematical Problems in Engineering* 2016, 1–9 (2016)
5. do Nascimento Vale, L., de A. Maia, M.: Keece: Mining key architecturally relevant classes using dynamic analysis. In: *2015 IEEE International Conference on Software Maintenance and Evolution (ICSME)*. pp. 566–570 (2015)
6. Fawcett, T.: An introduction to roc analysis. *Pattern Recognition Letters* 27(8), 861–874 (2006), <https://www.sciencedirect.com/science/article/pii/S016786550500303X>, rOC Analysis in Pattern Recognition

7. Gall, H., Hajek, K., Jazayeri, M.: Detection of logical coupling based on product release history. In: Proceedings of the International Conference on Software Maintenance. pp. 190–. ICSM '98, IEEE Computer Society, Washington, DC, USA (1998), <http://dl.acm.org/citation.cfm?id=850947.853338>
8. Kamran, M., Ali, M., Akbar, B.: Identification of core architecture classes for object-oriented software systems. *Journal of Applied Computer Science & Mathematics* 10, 21–25 (2016)
9. Meyer, P., Siy, H., Bhowmick, S.: Identifying important classes of large software systems through k-core decomposition. *Adv. Complex Syst.* 17 (2014)
10. Şora, I.: 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,. pp. 129–140. INSTICC, SciTePress (2015)
11. Şora, I., Chirila, C.B.: Finding key classes in object-oriented software systems by techniques based on static analysis. *Information and Software Technology* 116, 106176 (2019), <https://www.sciencedirect.com/science/article/pii/S0950584919301727>
12. Osman, M.H., Chaudron, M.R.V., v. d. Putten, P.: An analysis of machine learning algorithms for condensing reverse engineered class diagrams. In: 2013 IEEE International Conference on Software Maintenance. pp. 140–149 (2013)
13. Page, L., Brin, S., Motwani, R., Winograd, T.: The pagerank citation ranking: Bringing order to the web. Technical Report 1999-66, Stanford InfoLab (November 1999), <http://ilpubs.stanford.edu:8090/422/>, previous number = SIDL-WP-1999-0120
14. Pan, W., Song, B., Li, K., Zhang, K.: Identifying key classes in object-oriented software using generalized k-core decomposition. *Future Generation Computer Systems* 81, 188–202 (2018), <https://www.sciencedirect.com/science/article/pii/S0167739X17302492>
15. Şora, I.: 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). pp. 495–500 (May 2015)
16. Stana, A.D., Şora, I.: 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,. pp. 486–493. INSTICC, SciTePress (2019)
17. Steidl, D., Hummel, B., Juergens, E.: Using network analysis for recommendation of central software classes. In: 2012 19th Working Conference on Reverse Engineering. pp. 93–102 (2012)
18. Tahvildari, L., Kontogiannis, K.: Improving design quality using meta-pattern transformations: a metric-based approach. *J. Softw. Maintenance Res. Pract.* 16, 331–361 (2004)
19. Thung, F., Lo, D., Osman, M.H., Chaudron, M.R.V.: Condensing class diagrams by analyzing design and network metrics using optimistic classification. In: Proceedings of the 22nd International Conference on Program Comprehension. p. 110–121. ICPC 2014, Association for Computing Machinery, New York, NY, USA (2014), <https://doi.org/10.1145/2597008.2597157>
20. Yang, X., Lo, D., Xia, X., Sun, J.: Condensing class diagrams with minimal manual labeling cost. In: 2016 IEEE 40th Annual Computer Software and Applications Conference (COMPSAC). vol. 1, pp. 22–31 (2016)
21. Zaidman, A., Calders, T., Demeyer, S., Paredaens, J.: Applying webmining techniques to execution traces to support the program comprehension process. In: Ninth European Conference on Software Maintenance and Reengineering. pp. 134–142 (2005)
22. Zaidman, A., Demeyer, S.: 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)