



Identifying software components from object-oriented APIs based on dynamic analysis

Anas Shatnawi, Hudhaifa Shatnawi, Mohamed Aymen Saied, Zakarea Al-Shara, Houari Sahraoui, Abdelhak-Djamel Seriali

► To cite this version:

Anas Shatnawi, Hudhaifa Shatnawi, Mohamed Aymen Saied, Zakarea Al-Shara, Houari Sahraoui, et al.. Identifying software components from object-oriented APIs based on dynamic analysis. ICPC 2018 - 26th International Conference on Program Comprehension, May 2018, Gothenburg, Germany. pp.189-199, 10.1145/3196321.3196349 . lirmm-01932804

HAL Id: lirmm-01932804

<https://hal-lirmm.ccsd.cnrs.fr/lirmm-01932804>

Submitted on 11 Sep 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Identifying Software Components from Object-Oriented APIs Based on Dynamic Analysis

Anas Shatnawi
University of Milan-Bicocca
Milan, Italy
anas.shatnawi@unimib.it

Hudhaifa Shatnawi
Maharishi University of Management
Fairfield, Iowa, USA
hshatnawi@mum.edu

Mohamed Aymen Saied
Concordia University
Montreal, Quebec, Canada
med.aymen.saied@gmail.com

Zakarea Al Shara
University of Montpellier
Montpellier, France
zakarea.alshara@gmail.com

Houari Sahraoui
University of Montreal
Montreal, Quebec, Canada
sahraouh@iro.umontreal.ca

Abdelhak Seriai
University of Montpellier
Montpellier, France
seriai@lirmm.fr

ABSTRACT

The reuse at the component level is generally more effective than the one at the object-oriented class level. This is due to the granularity level where components expose their functionalities at an abstract level compared to the fine-grained object-oriented classes. Moreover, components clearly define their dependencies through their provided and required interfaces in an explicit way that facilitates the understanding of how to reuse these components. Therefore, several component identification approaches have been proposed to identify components based on the analysis object-oriented software applications. Nevertheless, most of the existing component identification approaches did not consider co-usage dependencies between API classes to identify classes/methods that can be reused to implement a specific scenario. In this paper, we propose an approach to identify reusable software components in object-oriented APIs, based on the interactions between client applications and the targeted API. As we are dealing with actual clients using the API, dynamic analysis allows to better capture the instances of API usage. Approaches using static analysis are usually limited by the difficulty of handling dynamic features such as polymorphism and class loading. We evaluate our approach by applying it to three Java APIs with eight client applications from the DaCapo benchmark. DaCapo provides a set of pre-defined usage scenarios. The results show that our component identification approach has a very high precision.

KEYWORDS

Software components, reverse engineering, object-oriented APIs, dynamic analysis, source code, understandability, reuse

ACM Reference Format:

Anas Shatnawi, Hudhaifa Shatnawi, Mohamed Aymen Saied, Zakarea Al Shara, Houari Sahraoui, and Abdelhak Seriai. 2018. Identifying Software

Components from Object-Oriented APIs Based on Dynamic Analysis. In *ICPC '18: ICPC '18: 26th IEEE/ACM International Conference on Program Comprehension*, May 27–28, 2018, Gothenburg, Sweden. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3196321.3196349>

1 INTRODUCTION

Decades of research have shown that developing software applications based on Application Programming Interfaces (APIs) improves software reuse by offering pre-implemented and tested functionalities [1] [2] [3] [4]. For Object-Oriented (OO) APIs, the basic unit is a class, which encapsulates its functionalities and specifies the public interface for using them [4]. Reusing and understanding large OO APIs, e.g., JDK and .NET framework APIs, are complex tasks due to large numbers of included classes and methods [1] [4] [5] [6]. Considering JDK 1.8.0 as an example, we have approximately 4240 classes used to provide its functionalities to software applications [7]. Thus, it is a challenge to understand this large number of classes to identify needed functionalities [8].

On the other hand, API classes/methods are reused based on reuse scenarios represented by the combinations of API classes or methods that offer the required functionalities to software applications [1] [4]. However, different applications rely on different scenarios, depending on their needs of API functionalities [8]. This leads to the possibility to have a large number of scenarios, corresponding to different combinations of API classes/methods [8]. To help support understanding these reuse scenarios, several research approaches have been proposed to abstract high level views of these scenarios in terms of frequent co-usage patterns between API classes/methods. The frequent co-usage patterns are identified based on the analysis of how software applications (called client applications) have (re)used API classes/methods [1] [9] [8] [10] [11] [12] [13]. These abstracted co-usage patterns of API classes/methods are used to support software engineers to perform several engineering tasks; API reengineering [1], API documenting [10], API understanding [11], API policy enforcing [14] [15], etc.

Identifying components in OO APIs is another alternative to improve APIs' documentation by providing an abstract high-level view of the provided functionalities. In addition to the documentation motivation, the identified components can also help migrating OO APIs into component-based ones if desired. Such a reengineering process does not only support feeding component-based repositories, but it also allows one to get the benefits of the component-based

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ICPC '18, May 27–28, 2018, Gothenburg, Sweden

© 2018 Association for Computing Machinery.

ACM ISBN 978-1-4503-5714-2/18/05...\$15.00

<https://doi.org/10.1145/3196321.3196349>

software engineering by developing large-scale software applications using selected components and integrating them together based on flexible architectures [16]. It has been admitted that software components are more reusable and understandable software modules than object-oriented classes. This admission is based on the granularity level where components are coarse-grained modules compared to the fine-grained object-oriented classes [17]. Furthermore, components clearly define their dependencies through their provided and required interfaces in an explicit way that facilitates the understanding of how to reuse these components. These are the motivations of several component identification approaches that have been proposed to identify components based on the analysis of source code of OO software [18] [16] [19].

Existing component identification approaches are designed to identify components from OO source code of software applications. They only rely on OO dependencies (e.g., method invocations, sharing types, etc.) to identify dependencies between classes. However, in the context of software OO APIs, co-usage dependencies between API classes should also be considered to identify classes/methods that can be reused to implement a specific scenario. This requires the analysis of source code of APIs and their client applications. Therefore, approaches designed for OO applications cannot be applied directly to OO APIs. To the best of our knowledge, the only approach proposed to identify components from OO APIs is one described in [1] and [20]. This approach is based on the static analysis of the source code of OO APIs and their client applications. As we are dealing with actual clients using the OO API, dynamic analysis allows to better capture the instances of API usage. Approaches using static analysis are usually limited by the difficulty of handling dynamic features such as polymorphism and class loading.

In this paper, we propose an approach to identify reusable components in OO APIs. Our approach is based on the dynamic analysis of interactions between client applications and the targeted OO API. Thus, we generate the execution traces for the different usage scenarios implemented in client applications. These execution traces realize dependencies between API classes through their method invocations. We assume that components are identified in terms of API classes that are frequently reused together by client applications. Therefore, we consider that groups of methods frequently appearing together in execution traces form provided interfaces of components, and the owner API classes of these methods constitute the structure of that candidate component. To evaluate our approach, we applied it to three Java APIs with their clients from the DaCapo benchmark. The results show that the precision of our component identification approach is 98%.

The rest of this paper is organized as follows. Section 2 presents the framework behind our approach. Section 3 describes the process of identifying execution traces related to usage scenarios of client applications of APIs. In Section 4, we identify graph representations of APIs methods based on the relationships appeared in execution scenarios using a proposed quality function. Section 5 shows how classes composing components are identified based on a graph-based clustering algorithm. Evaluation results are presented in Section 6. Related works are discussed in Section 7. Conclusion and future work are presented in Section 8.

2 THE PROPOSED APPROACH FRAMEWORK

In this section, we provide an overview, summarize the principles and propose the process of the proposed approach.

2.1 Approach Overview

Our approach identifies reusable components based on the dynamic analysis of interactions between client applications and the targeted API. We define a software component as a collection of classes that participate to implement one or more functionalities for client applications of an API. The provided interfaces, of this component, are API methods that have been invoked together frequently by client applications. Conversely, the required interfaces are API methods that have been invoked by the component classes and belong to other API components' classes.

Classes composing a component are identified based on their invoked methods (provided interfaces). To determine which methods are invoked jointly, we analyze execution traces of client applications using the API. As for dynamic analysis, we need representative usage scenarios, we rely on use cases of client applications to identify such representative scenarios, and execute them to produce the traces. The execution traces highlight the dependencies between API methods. We consider methods that appear frequently together in execution traces as provided interfaces of a component. The owner classes of these methods define the structure of the component implementation.

We want to package each collection of API methods frequently used together to form a provided interface of a reusable component, without changing the internal structure of the API (i.e., we do not aim to identify architectural view of the API). Therefore, we allow a class to be a part of more than one component as different subsets of its methods can participate with various groups of methods related to other classes to implement different functionalities.

To evaluate the quality of a candidate collection of API methods to form a component provided interface, we define quality function that analyzes dependencies between API methods based on the identified execution traces. We cluster API methods using this quality function.

2.2 Approach Principles

The principles of our approach are summarized as follows.

- A component is defined as a group of classes identified through its potential provided interfaces.
- A provided interface of a component is a set of methods used frequently together in execution traces.
- A required interface of a component is a set of methods used by this component and belonging to other components' classes.
- A class can belong to several components since different subsets of its methods can be included in different component interfaces.
- Execution traces are used to identify the dependencies between API methods. Thus, they guide the component identification process.

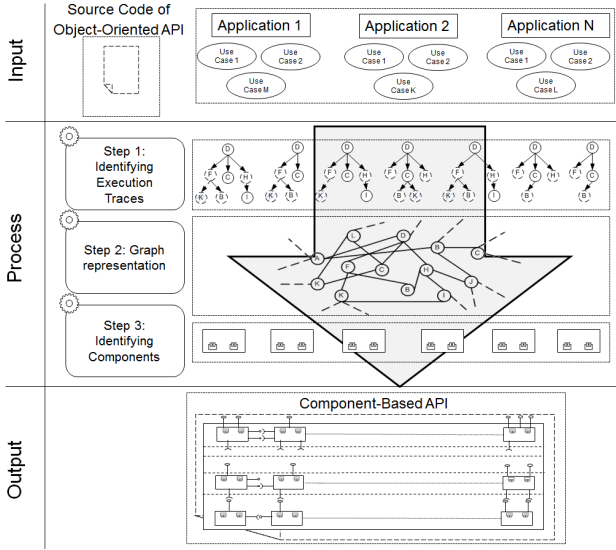


Figure 1: The process of component identification from object-oriented APIs

2.3 Approach Process

To implement our approach, we propose a process, presented in Figure 1, based on the following three steps:

- (1) **Identifying execution traces:** execution traces are identified based on usage scenarios related to client applications of APIs. Each executing trace is realized in terms of a call tree that represents the dynamic relationships between API methods corresponding to a usage scenario.
- (2) **Building graph representations of APIs:** to facilitate the problem definition, we rely on graph representations to highlight dependencies between methods of an API. These dependencies are identified based on the strength of relationships between methods in the call trees.
- (3) **Identifying classes composing components:** we apply a graph-based clustering algorithm to partition the graph into sub-graphs. Each sub-graph represents a group of methods forming a provided interface of a component, while their classes represent the component implementation.

3 IDENTIFYING EXECUTION TRACES

In this section, we discuss how to identify execution traces related to APIs based on usage scenarios of client applications.

An execution trace is a set of methods that represent the execution of a usage scenario of a client application of the targeted API. These methods can be represented in terms of a call tree defined as a *directed tree* $T = (V, E)$, where V is a set of vertices referring to a set of API methods and E is a set of edges. An edge $(V1, V2)$ refers to that a method $V1$ invokes a method $V2$. The root of the tree represents the starting point of the execution trace.

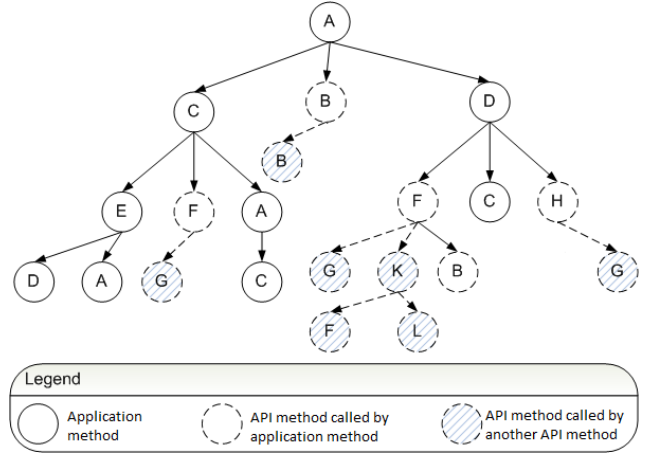


Figure 2: Example of call tree corresponding to usage scenario

3.1 Executing Usage Scenarios to Identify Call Trees

We execute each usage scenario to identify the set of methods corresponding to its implementation. A method representing the entry of a usage scenario is considered as the call tree root. Then, we follow method calls to identify the other nodes. When a method invokes a method, a new node is added to the tree with an edge from the first node to the class node, and so on.

Methods related to the call tree's nodes can be classified into three categories: methods implemented in application's classes, method of API classes used by application's classes to access API functionalities and method of API classes that have got called by the second category's methods (from API method used by method of application's classes). Figure 2 shows an example of a call tree that is composed of 4 application methods, which are A, C, D and E, and 6 API methods, which are B, F, G, H, L and K. The B and F methods have been invoked from an application class, while G and K methods have been called from an API method.

3.2 Removing Application's Methods from Call Trees

Since our goal is to analyze interdependencies between API methods, we prune the identified call trees by removing methods related to application classes. For each call tree, we perform a *Breadth First Search* algorithm starting from the root node. If a child is an API method, then we leave it. Otherwise, we remove the node. The children of the removed node are attached to the parent of the removed node. This process is continued until verifying all tree levels. For example, Figure 3 shows the pruned tree of the call tree presented in Figure 2. As it is noticed, the tree size is significantly reduced which decreases the complexity of analyzing call trees. For instance, it is reduced from 21 to 13 nodes.

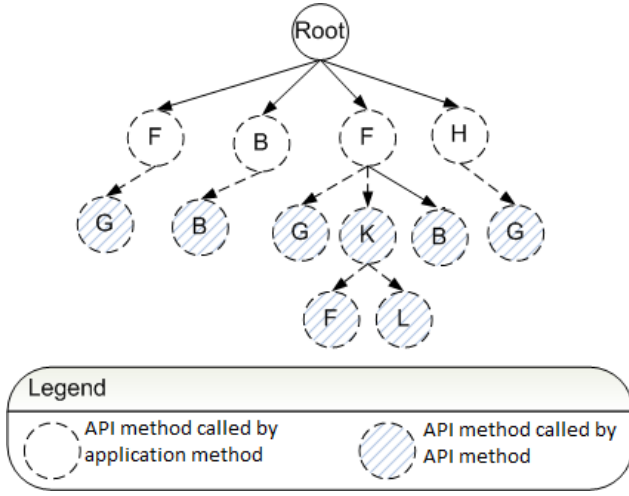


Figure 3: Pruned call tree of one in Figure 2

4 BUILDING GRAPH REPRESENTATIONS OF APIS

We represent an API in terms of an *undirected weighted graph* $G = \langle V, E \rangle$, where V is the set of API methods and E is the set of edges between the vertices V . An edge links the vertices (u, v) if they appeared together in call trees. The weight associated to each edge is based on the strength of relationship of its vertices in call trees. We distinguish three attributes that should be analyzed in call trees to measure the strength of these relationships. These are call frequency, call distance and call weight.

4.1 Call Frequency

This attribute refers to the relationship between numbers of co-occurrence of a group of methods in call trees and their cohesive. Methods frequently appeared in call trees are likely participating to provide related functionalities. For instance, in the call trees presented in Figure 4, the call frequency of A and B is 2 times, and it is 1 for A and L . Thus, A and B are likely cohesive than A and L .

The call trees are identified from different applications developed by different development teams. Thus, we propose to measure the call frequency based on two cases related to methods frequently used together in call trees of the same applications by the same team or different applications of different teams. We consider that methods used by different development teams will provide a global view of their reuse frequency (i.e., global frequency) compared to methods used by the same development team (i.e., local frequency).

Local frequency measures how many methods are used together in the same applications (the same developers). It is measured based on the average number of appearances of a group of methods in one application. We calculate the ratio between the number of call trees containing the methods to the total number of call trees in each application (c.f. Equation 2).

Global frequency measures how many methods are used together in different applications. Global frequency is calculated based on the average number of applications that contain the methods to the

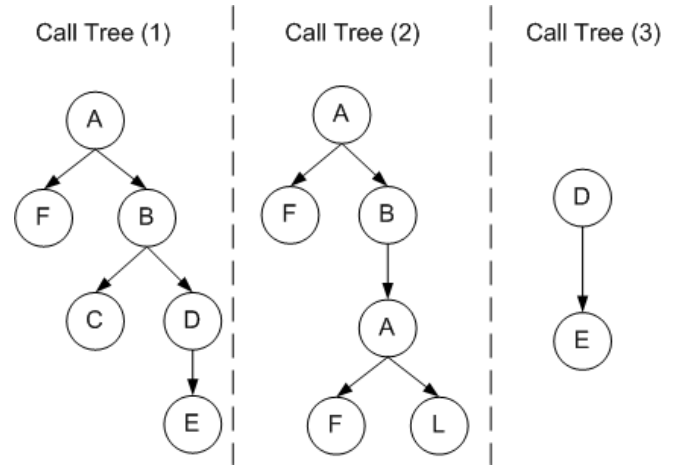


Figure 4: An example of three call trees

total number of applications (c.f. Equation 3). The call frequency of a group of methods is the mean value of local and global frequencies of all pairs of methods.

Equation 1 measures the call frequency for a set of methods E , where $Co-occur(c, v)$ returns 1 if the methods c and v exist in the call tree t , otherwise it returns 0. $T(a)$ refers to the set of call trees of the application a in the set of applications $apps$.

$$CallFreq(E) = \frac{\sum_{c, v \in E, c \neq v} \frac{LFreq(c, v) + GFreq(c, v)}{2}}{(|E| * (|E| - 1)) / 2} \quad (1)$$

$$LFreq(c, v) = \frac{\sum_{a \in apps} \frac{\sum_{t \in T(a)} Co-occur(c, v)}{|T(a)|}}{|apps|} \quad (2)$$

$$GFreq(c, v) = \frac{NumberOfAppsContaining(c, v)}{|apps|} \quad (3)$$

4.2 Call Distance

This attribute is related to the *distance* between methods in call trees, which refers to the strength of their interactions. As much as methods are closer in call trees they have high probabilities to provide same functionalities. For example, in the first call tree in Figure 4, the distance between A and B is 1, while it is 3 between A and E . Thus, A and B have higher probability to be cohesive than A and E .

We define Equation 4 to measure the call distance of a group of methods E . For each pair of methods, it calculates the mean distance of this pair in all call trees using Equation 5. Then, the call distance is the average value of the mean values of all pairs.

$$CallDist(E) = \frac{\sum_{c, v \in E, c \neq v} Distance(c, v)}{(|E| * (|E| - 1)) / 2} \quad (4)$$

$$Distance(c, v) = \frac{\sum_{a \in apps} \frac{\sum_{t \in T(a)} dis(c, v, t)}{|T(a)|}}{|apps|} \quad (5)$$

$$dis(c, v, t) = 1 - \left(\frac{avgDistance(c, v, t)}{2 * getTreeDepth(t)} \right) \quad (6)$$

Where $avgDistance(c, v, t)$ finds the average distance between the methods c and v in a call tree t . It returns $2^*getTreeDepth(t)$ if the methods c and v do not appear in a call tree.

4.3 Call Weight

This attribute is related to the weight of the co-occurrence of a group of methods in call trees. Methods that are co-located in the same call trees have different weight of co-occurrences. This affects to the degree of their interaction, and consequently their cohesive. The weight depends on the number of calls in the call trees. For example, in Figure 4, D and E appear together in call tree (1) and call tree (3), but with different weights. This means that their appearance has different impact to their cohesive. In fact, it is clear that their appearance in call tree (3) has a higher impact than in call tree (1).

We measure the call weight based on the number of edges in call trees. For a group of methods, the call weight is the average call weight of all pairs in this group. Since the call weight of a pair of methods in such a call tree is the fraction between the number of direct call between these methods and the total number of edges in the call tree. For our example (D and E), it is 20% in call tree (1) and 100% in call tree (3), thus it is 60% in average. Equation 7 measures the call weight of a collection of methods E :

$$CallWeight(E) = \frac{\sum_{c, v \in E, c \neq v} Weight(c, v)}{(|E| * (|E| - 1)) / 2} \quad (7)$$

$$Weight(c, v) = \frac{\sum_{a \in apps} \sum_{t \in T(a)} Wei(c, v, t)}{|apps|} \quad (8)$$

Where $Wei(c, v, t)$ is the wight between the methods c and v in the call tree t . It is calculated using Equation 9.

$$Wei(c, v, t) = \frac{NumberOfDirectCall(c, v, t)}{TotalNumberOfCall(t)} \quad (9)$$

For evaluating the quality of a group of methods E , we define Equation 10 as a linear combination of the three functions evaluating each attribute.

$$Q(E) = \frac{1}{\sum_i \lambda_i} \cdot (\lambda_1 \cdot CallFreq(E) + \lambda_2 \cdot CallDist(E) + \lambda_3 \cdot CallWght(E)) \quad (10)$$

Where λ_1 , λ_2 and λ_3 are weight values, situated in $[0-1]$. These are used by the API expert to weight each attribute as needed.

5 IDENTIFYING COMPONENTS THROUGH THEIR PROVIDED INTERFACES

As we mentioned previously, methods composing provided interfaces of components are grouped based on a graph-based clustering algorithm. We select a clustering algorithm that allows overlapping between clusters. The idea behind that is to allow a method to appear in different clusters (component provided interfaces) as when it is used frequently with different groups of method.

5.1 Identifying Overlapping Clustering of Methods

The identification of methods composing component interfaces is based on the identification of a set of subgraphs, where they have accepted quality values. We define the component provided interface

as a subgraph $G^* = \langle V^*, E^* \rangle$ where the average weight of edges in E^* maximizes the quality function, the set of methods corresponding to the vertices V^* represents the component provided interfaces. The identification of an optimal set of subgraphs needs to identify all possible candidate subgraphs that can be extracted from the graph. Then, the selection of subgraphs is made by choosing ones that maximize the quality function. Nevertheless, this is considered as *NP-Complete* problem since the time complexity is exponential. Thus, we propose to use a heuristic clustering algorithm to find a near optimal set of subgraphs. We use the *OClustR* [21] clustering algorithm. We select this algorithm as it is a graph-based clustering algorithm that is able to identify overlapping clusters. In addition, there is no need to identify the number of the needed clusters as we do not have any idea about the number of components that should be identified.

Based on *OClustR*, the problem is to identify a set of subgraphs $W = \{G_1^*, G_2^*, \dots, G_k^*\}$, where the graph G is covered by W (i.e. $\bigcup_{i=1}^k V_i^* = V$, such that G_i^* covers a vertex $v \in V$ if v belongs to V_i^*). Finding the minimum number of subgraphs, i.e. the set W , is known as a vertex cover problem. A vertex cover is a smallest subset of vertices $V^\wedge \subseteq V$, such that each vertex $v \in V$ is either a part of V^\wedge or directly connected, i.e. an adjacent vertex, to another vertex $u \in V^\wedge$. Nevertheless, a vertex cover is classified as *NP-Complete* problem [Ref]. Thus, *OClustR* proposed an approximation algorithm that aims at mining a near optimal set of vertices V^\wedge , such that graphs corresponding to V^\wedge cover G and maximizing the quality of the corresponding component as well. This algorithm identifies the subgraphs based on two steps. The first step aims at mining the initial subgraphs (clusters). The second step goals at refining the initial subgraphs by minimizing the number of clusters, since they may include useless subgraphs, and the overlapping between them.

5.1.1 Mining initial clusters. This step aims at identifying the set of initial subgraphs W , such that each subgraph G_i^* is a *weighted star subgraph*, (*ws-graph*) for short, in G . A *ws-graph* $G_i^* = \langle V^*, E^* \rangle$ have a vertex $c \in V^*$, c is called the *center* of G_i^* , where there is an edge connecting c with the other vertices in V^* , the other vertices are called *satellites*. Here, the problem is to identify a set of *center* vertices $Z = \{c_1, c_2, \dots, c_k\}$, where each vertex $c_i \in Z$ is the center of the subgraph $G_i^* \in W$. In order to identify the set Z , all vertices should be investigated. The investigation is based on an iterative process. At each iteration, one vertex v is added to Z . This is continued until Z reaches the condition of covering G . The selection of a vertex v to be added to Z is based on an evaluation criterion. This criterion depends on two factors. The first one aims at maximizing the cover of G at each iteration and controlling the overlapping between the *ws-graphs*. This factor is called *Relative Density*, (*RD*) for short. The second factor goals at keeping the quality of the component corresponding to the *ws-graph* related to a vertex v , since a vertex that provides a high coverage of G could produce a low quality component corresponding to our quality fitness function. This factor is called *Relative Compactness*, (*RC*) for short.

RD is related to the number of *satellite* vertices that could be covered by *ws-graph* corresponding to a *center* vertex v . Since vertices having higher degree, i.e. vertices having more *satellites*, produce *ws-graphs* that contribute more to cover G , the selection of a vertex

should take the vertex degree into account. In some cases, this way is not sufficient since some of the *satellite* vertices maybe already covered by another *center* vertex r that is added to Z in a previous iteration. Thus, we need to take the number of *satellite* vertices that will be covered by selecting a vertex v in the current iteration, i.e. we exclude *satellite* vertices already covered by another *center* vertices. We measure RD of a vertex $v \in V$ based on Equation 11. The higher value of RD the better vertex to be selected, where its value is situated in $[0-1]$.

$$RD(v) = \frac{\text{NumberOfUncoverSatelliteVertices}}{\text{TotalNumberOfSatelliteVertices}} \quad (11)$$

RC is related to the quality of a *ws-graph* corresponding to v compared to its *satellite* vertices. This means that we investigate if a vertex v is the best vertex that maximizing the quality among its *satellite* vertices. Therefore, we measure RC based on the ratio between the number of *satellite* vertices that produces *ws-graphs* having lower quality values than the *ws-graph* corresponding to v and the total number of *satellite* vertices. Where the quality of a *ws-graph* is the average weight values between all pairs of vertices included by the *ws-graph*. We calculate RC based on Equation 12, where $SatelliteCompactness$ is the number of *satellites* s such that the quality of *ws-graph* of s is grater than the quality of the *ws-graphs* of v . The higher value of RC the better vertex to be selected, where its value is situated in $[0-1]$.

$$RC(v) = \frac{\text{SatelliteCompactness}}{\text{TotalNumberOfSatelliteVertices}} \quad (12)$$

Based on these factors, we firstly sort the vertices in a decreasing order based on their relative quality; the average of their RD and RC , denoted by RQ . Then, these vertices are iteratively added to the set Z with respect to one of these conditions: (1) a vertex is not covered and (2) it is covered but at lease there is one of its *satellites* that is not covered. This continues until covering graph G .

Algorithm 1: Identifying initial clusters

Input: undirected weighted graph $G = \langle V, E \rangle$
Output: A Set of center vertices Z
 $L \leftarrow G(V)$;
Sort L in a decreasing order based on RQ ;
for each $v \in L$ **do**
 if $v.uncovered() \parallel v.hasUncoveredsatellites()$ **then**
 $Z \leftarrow Z \cup \{v\}$;
 end
end
return Z ;

5.1.2 Refining the initial clusters. This step aims at enhancing the initial clusters, identified in the previous step. The goal is at reducing the number of resulted clusters as well as the overlapping between them. The process of identifying the set Z is a greedy one since it selects the vertex having the highest quality at each iteration. Thus, this may lead to the situation of adding a useless vertex u to Z . The identification of a useless subgraph G_u^* is based on how much its *satellites* are shared with other G^* s. The worst case is that $u \in Z$, such that G is still covered by the G^* s corresponding

to $Z - \{u\}$, i.e. u and its *satellites* are covered by other G^* s. The average case is that $u \in Z$, such that G_u^* shares most of its *satellites* with the other G^* s corresponding to $Z - \{u\}$. In these case, the G_u^* is considered as a useless subgraph and u should be deleted from Z .

According to that, G_u^* is considered as a useless if it meets two conditions. The first one is that u is a *satellite* vertex of at lease another G_v^* . The second condition is that V_u^* shares more than half of its *satellite* vertices with another G^* s. Once a useless G_u^* is identified, u is removed from Z , and the non-shared *satellites* are distributed to another G^* . We select the G^* that covers u as it covers the center of G_u^* . In case that there is more than one G^* covering u , we select G^* having the grater number of satilletes. By doing this, we allow producing clusters having many vertices.

Algorithmically, we refine the initial clusters G^* as follows. It firstly sorts the vertices of Z in a descending order based on their degrees and sign all of them as unvisited. Then, starting from the vertex having the highest degree, each vertex v is analyzed. The analysis consists of removing from Z any vertex $u \in Z$ such that u is a *satellite* of G_v^* ($u \in V_v^*$) and G_u^* is considered as useless G^* . Then, all *satellite* verices of G_u^* is added to G_v^* since G_v^* is the one having the grater number of satilletes among u adjacent vertices. Once all *satellite* verices in G_v^* are analyzed, v is signed as visited and G_v^* is considered as a final cluster.

Algorithm 2: Refining initial clusters

Input: A Set of center vertices Z
Output: A Set of Clusters C
 $L \leftarrow G(V)$;
Sort Z in a decreasing order based on their degree;
for each $v \in Z$ **do**
 $v.isVisited(false)$
end
for each $v \in Z$ **do**
 for each $u \in \text{satellites of } G_v^*$ **do**
 if $u \in Z \ \&\& \ u.isNotVisited()$ **then**
 if G_u^* is useless **then**
 $G_v^* \leftarrow G_v^* \cup G_u^*$;
 $Z \leftarrow Z - \{u\}$;
 end
 else
 $u.isVisited(true)$
 end
 end
 end
 $C \leftarrow C \cup \{V_v^*\}$;
end
return C ;

6 EVALUATION RESULTS

6.1 Description of APIs and their Client Applications

We evaluate our approach using real usage scenarios from DaCapo Benchmarks [22]. We select DaCapo because it is composed of a

Table 1: Data set description

API name	Size (classes)	# of app. clients	App. client names
ACL	23	2	Tomcat, Fop
ASM	99	2	Jython, Pmb
XML	302	4	Xalan, Pmd, Fop, Batik

set of open-source real Java applications coupled with collections of *pre-defined usage scenarios*.

To show the applicability of our approach in several context, we select from DaCapo three widely used Java APIs of three different size (i.e., 23, 99 and 302 classes) and from three different domains. The first one is the Apache Commons Logging (ACL) API¹ that contains 23 classes, and offers a log interfaces and a middleware/tooling developer with a simple logging abstraction. The second one is the ASM API² consists of 99 classes. It provides functionalities related to Java bytecode manipulation and analysis. The third API is the XML API³ that is composed of 302 classes. It offers functionalities for processing XML files.

DaCapo only provides a limited number of usage scenarios of software applications developed based on these three APIs. We are able to identify only eight client applications that use methods of classes of the APIs. These applications are as follows.

- **Tomcat** implements J2EE technologies like Servlet and JavaServer Pages.
- **Fop** is an output-independent print formatter that parses and formats XSL-FO files.
- **Xalan** is an XSLT processor for transforming XML documents.
- **Batik** is a Scalable Vector Graphics (SVG) toolkit that renders a number of SVG files. Xalan, Fop and Batik use classes from the XML API.
- **Pmd** is a source code analyzer for Java code. It uses two APIs; XML and ASM ones.
- **Jython** is a python interpreter written in Java to execute and interpret Python programs. It is a client for the ASM API.
- **Lusearch** and **Luindex** are respectively a text search tool and a text indexing for a corpus of data comprising the works of Shakespeare and the King James bible. Both applications are clients of the Lucene API.

Table 1 shows the description of the collected dataset. For each API, we provide the number of included classes, the number of applications considered as clients of this API and the names corresponding to the client applications.

6.2 Evaluation Process

Our evaluation process is based on four steps. First, we have discussed the results of identifying call trees related to execution traces of usage scenarios of client applications of APIs. Then, we have shown the results of identifying groups of methods that are considered as provided interfaces of components. Next, we have presented

¹Available at: <https://commons.apache.org/proper/commons-logging/guide.html>

²Available at: <http://www-etud.iro.umontreal.ca/saiedmoh/asm-3.3.1/doc/javadoc/user/index.html>

³Available at: <https://xerces.apache.org/xerces2-j/javadocs/api/index.html>

how we evaluate the resulting components based on functionalities provided by methods representing their provided interfaces. Next, we have provided a discussion about dynamic and static analysis component identification approaches. Finally, we have discussed threats to validity related to the results of our approach.

6.3 Results of Identifying Call Trees Based on Executing Usage Scenarios

We rely on BTrace⁴ dynamic tracing tool to collect execution traces related to usage scenarios. BTrace supports the execution of a Java program based on its bytecode and dynamically collect the methods related to this execution. We configure BTrace to run each client application only once since we are not interested in performance analysis.

Table 2 shows the results related to the identified call trees of execution traces of usage scenarios of client applications of APIs. For each call tree, we present the tree's size in terms of the number of included nodes, the number of distinct API methods included in this call tree, the tree's height, the minimum, maximum and average method' repetition in this call tree.

The results show that the call trees have small heights compared to the large number of nodes included in the trees. The average height of call trees related to ACL, XML and ASM APIs is respectively 12.5 $((12+13)/2)$, 7.5 $((7+4+14+5)/4)$ and 18.5 $((10+27)/2)$, while the average trees size is respectively 472676 and 1946555 nodes. This means that execution traces go deeply by 12.5, 7.5 and 18.5 methods in average respectively for ACL, XML and ASM APIs.

Indeed, the root of a call tree represents an execution scenario, and its children nodes represents the set of methods directly invoked in the source code of client applications to access API's functionalities. Each sub-tree that corresponds to a root child explains internal dependencies related to method invocation between API methods to provide the invoked functionality by the corresponding root child. For example, the height of a call tree of 18 nodes refers to the maximum number of API methods that are invoked corresponding to an API method invocation directly invoked by client applications.

Another observation is that the call trees are widely distributed in a horizontal way, i.e., each node in a call tree has a large number of children relatively to its height. Therefore, API methods have dense interactions with each others such that each API method invokes a bunch of other API methods. In addition to that, we find that some API methods have been invoked only once compared to some other ones which have been invoked 198972.5 times in average (average max node repetition). Many groups of methods are invoked in similar frequently at the same time. These are interested to be in a same provided interface, even if they belong to different classes.

6.4 Results of Component Identification

The results of our clustering algorithm are shown in Table 3. For each API, we present the number of identified components, the average size of provided interfaces in terms of methods and the average component size in terms of included classes.

Respectively for ACL, ASM and XML APIs, the results show that the average number of methods that need to be used together to

⁴Available at: <https://kenai.com/projects/btrace>

Table 2: The results of call trees identification

Call tree	Total no. of nodes	No. of API unique methods	Tree height	Min method repetition	Max method repetition	Average method repetition
ACL1	1983448	40	12	1	971877	49586.2
ACL2	467328	81	13	1	233231	5769.48
XML1	79844	57	4	1	11040	1400.77
XML2	116013	79	7	1	25392	1468.52
XML3	245	25	14	1	20	9.8
XML4	1694602	84	5	1	769600	20173.83
ASM1	3366089	106	10	1	589875	31755.56
ASM2	527021	28	27	1	232285	18822.18

Table 3: The results of component identification

API	# of identified components	Avg. provided interface size (methods)	Avg. component size (classes)
ACL	9	11	2.14
ASM	24	5.08	2.13
XML	22	6.72	3.25

Table 4: Examples of identified components

API	Component provided interface names
ACL	org.apache.commons.logging.impl.SimpleLogrun, org.apache.commons.logging.impl.LogFactoryImplcreateLogFromClass, org.apache.commons.logging.impl.SimpleLoggetContextClassLoader, org.apache.commons.logging.impl.SimpleLog, org.apache.commons.logging.impl.SimpleLogclass, org.apache.commons.logging.impl.SimpleLogaccess, org.apache.commons.logging.impl.SimpleLoggetResourceAsStream, org.apache.commons.logging.impl.LogFactoryImpldiscoverLogImplementation, org.apache.commons.logging.impl.LogFactoryImplnewInstance, org.apache.commons.logging.impl.SimpleLog
ASM	org.objectweb.asm.commons.VisitorvisitCode, org.objectweb.asm.commons.VisitorvisitTableSwitchInsn, org.objectweb.asm.commons.VisitorvisitIntInsn, org.objectweb.asm.commons.Visitorvisit, org.objectweb.asm.commons.VisitorvisitMethodInsn, org.objectweb.asm.commons.VisitorvisitInsn, org.objectweb.asm.commons.VisitorvisitJumpInsn, org.objectweb.asm.commons.VisitorvisitMethod
XML	javax.xml.parsers.DocumentBuilderFactoryisIgnoringElementContentWhitespace, javax.xml.parsers.DocumentBuilderFactoryisNamespaceAware, javax.xml.parsers.DocumentBuilderFactory, javax.xml.parsers.DocumentBuilderFactoryisCoalescing, javax.xml.parsers.DocumentBuilderFactoryisExpandEntityReferences, org.xml.sax.InputSourcegetInputStream, javax.xml.parsers.SecuritySupportgetContextClassLoader, javax.xml.parsers.DocumentBuilder, javax.xml.parsers.DocumentBuilderFactorynewInstance, javax.xml.parsers.DocumentBuilderparse, javax.xml.parsers.SecuritySupport, org.xml.sax.InputSourcegetCharacterStream, javax.xml.parsers.DocumentBuilderFactoryisValidating, javax.xml.parsers.DocumentBuilderFactoryisIgnoringComments, org.xml.sax.InputSourcegetEncoding

access an API functionality is 9, 24 and 22 methods. These methods offer functionalities of 2.14, 2.13 and 3.25 owner classes in average.

For each API, Table 4 shows an example of a group of object-oriented API classes representing the implementation of identified component provided interfaces. To enable their reuse for software developers, these groups need to be transformed to confirm an existing component model. It worths to mention some works that can be used by the user of our approach to transform these object-oriented

implementation of identified components to be confirmed to component models. Alshara et al. [23] [24] provided initial approaches for transforming object-oriented component implementation to be confirmed to several component models (e.g., OSGi [25] and Fractal [26]).

6.5 Results of Evaluating Identified Components

The identified components are evaluated based on the functionalities offered by their provided interfaces. For each component, we evaluate how much of the methods composing its provided interface are related to provide the same API functionalities. The quality of each provided interface is calculated based on the ratio between the number of related methods to the total number of methods composing this interface (c.f. Equation 13). We rely on the API documentations to identify each method's functionality. Then, among a group of methods composing the provided interface, we select ones that have related functionalities to at least one method in this group.

For a pair of methods, we consider them as functional-related if one of these two cases is applied. The first one refers to identify direct indication(s) in the API documentations stated the correlation. The second one is based on human experts where their experiences allow them to decide if the two methods is functional-related. The number of selected methods represents the numerator in Equation 13. The authors perform the evaluation themselves. To avoid bias, each component was evaluated by at least three authors. The average represents the final evaluation result of this component.

$$Evaluation(Component) = \frac{NoOfRelatedMethods}{TotalNo.OfMethods} \quad (13)$$

The precision of our approach is the average of the resulting values by applying Equation 13 to components of an API. Figure 5 shows the precision. The results show that our approach is able to correctly identify components where their provided interfaces are composed of functional-related API methods with 98% precision $((100\%+93\%+100\%)/3)$. During our evaluation, we note that a component can be composed of one class and a subset of its methods forms the provided interface of this component. The same class forms another component with different subset of methods as provided interface. We consider these as correctly identified components that offer (sub)functionalities related to their provided interfaces. Each component can be (re)used as stand-alone or composed with other ones to generate higher level components (composite components).

Moreover, we observe that the same set of classes can form different components using different subsets of their methods as provided interfaces, i.e., component having the same classes but different provided interfaces. By checking the functionality of each method in these components, we find that they are related to same functionalities, but in different contexts. As it is noticed, some classes are parts of several components since different subset of their methods are part of these components. We look to the position of these classes as correct when the corresponding methods are harmonic with other methods in a given component. In XML API, we identify a component that its provided interface consists of 14 methods belonging to 11 different classes. After checking the functionality offered by these methods, we find that methods of 7 classes are related to the same functionality, while the other 4 classes are not related to each other since they are helper classes.

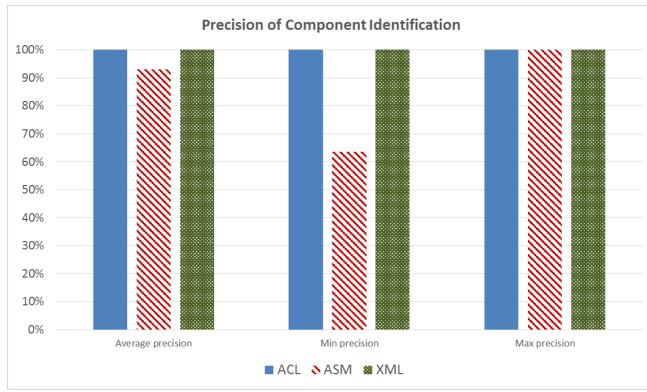


Figure 5: The precision of component identification

6.6 Dynamic vs Static Analysis

The dynamic analysis approach is very expensive (run-time startup, collect usage scenarios ...), if we compare it to the static analysis which can only rely on source code files. However, dynamic analysis is more accurate as it contributes to provide real dependencies that cannot be recovered only by static analysis. Such dependencies are those related to late-dynamic binding, method overriding polymorphism and implicit dependencies in configuration files. These dependencies can exist extensively or rarely depending on the application context and framework.

On the other hand, the dynamic analysis approach relies on usage scenarios to cover API functionalities. However, some API functionalities may not be covered by the usage scenarios implemented in the selected client applications (thus, they will not be covered by the identified components). We recommend to select client applications that cover the maximum number of possible usage scenarios, which can be a strong constraint in some cases.

We compare the two approaches by applying our approach and the static approach presented in [1] and [20] into the collected APIs and client applications. We find that both approaches construct different abstractions of the same functionalities. Components identified by dynamic analysis can be composed to form components resulting from static analysis. In terms of functionality-related, the

two approaches provide similar results. We interpret this similarity by the fact that the studied APIs do not rely on the dependencies that are not detected based on the static analysis approach (e.g., implicit dependencies). In this context, we would like to recommend to rely on the dynamic analysis approach for APIs that do rely on late-dynamic binding, methods overriding polymorphism and implicit dependencies in configuration files to codify dependencies.

6.7 Threats to Validity

We identify two types of threats to validity concerning the proposed approach; internal validity and external validity.

6.7.1 Threats to Internal Validity. Two aspects concern the internal threats to validity:

- (1) We need to rely on human expert opinions to evaluate the results. However, it is not easy to access the experts of the APIs. The authors considered themselves as experts since each author has at least more than 8 years of experience in software development.
- (2) Our evaluation does not consider the recall as we are not able to identify how many API methods are actually related to a functionality and not extracted by the approach. However, the recall depends on the coverage of the usage scenarios and does not rely on our approach. Any API method used in at least one usage scenario will be a part of at least one component.

6.7.2 Threats to External Validity. We identified two aspects related to the external threats to validity:

- (1) We evaluated our approach based on APIs and client applications written using *Java*. The proposed approach can be generalized not only for other object-oriented languages (e.g., C++), but also for procedural ones (e.g., C, Pascal). The idea is that our approach does not care about the implementing programming language since its input is a set of call trees. The nodes in these call trees can refer either to methods or procedures/routines. The main difference is that identified components will be groups of procedures.
- (2) The selection of API client applications may impact identified components as different applications use API methods following different scenarios. This may impact the reusability of identified components for new independent applications. However, we assume API functionalities are (re)used following similar patterns by developers. This assumption was successfully utilized and proven in several research papers [1] [8] [10] [12]. We recommend selecting as much as possible of API client applications to minimize the influence of domain-specific API usages.

7 RELATED WORK

To the best of our knowledge, there is only one approach proposed to identify components from object-oriented APIs [1] [20]. It identifies components as groups of API classes that are frequently used together by the client applications of an API and structurally dependent. To identify co-usage relationships between classes, Frequent-Pattern Growth algorithm was used where transactions are collections of API classes used by a client application. However, the

approach presented in [1] [20] does not trace the API method calls from client applications to deepen in the API as the proposed approach does to identify relationships between API classes. Instead, it relies on a static analysis technique to analyze structural dependencies in the source code of APIs. Recently developed software systems contains implicit dependencies resulted from the use of late dynamic-binding, Java-reflection and container-services offered by the frameworks. These implicit dependencies are difficult to be detected using static analysis techniques [27].

All of these existing component identification approaches that are designed to identify components from OO software applications are not be directly applied to OO APIs. The reason is that they only relied on OO dependencies (e.g., method invocations, sharing types, etc.) to identify dependencies between classes, while co-usage dependencies between API classes should also be considered to identify classes/methods that can be reused to implement a specific scenario.

In the following three paragraphs, we discuss the related work on component identification in standalone applications using dynamic and static analysis, as well as the identification of usage patterns.

Dynamic Analysis Component Identification from Applications. In [19], the authors presented an approach to identify component-based architectures from object-oriented applications. The identification is based on execution traces derived from use case scenarios. Classes frequently appearing in traces represent a candidate component. Grouping these classes is done by a clustering algorithm and a heuristic search, simulated annealing. Both techniques rely on a quality function. This quality function only considers the call frequency metric, which is insufficient without including metrics like our call weight and call distance. In [28], the authors presented an approach to extract an architecture view of object-oriented software using the analysis of execution traces. This research contribution lacks details on how such a process can be automated, and only the visualization aspect is discussed. In [29], the authors relied on the execution traces to recover two artifacts: 1) information that is used to update use cases' documentations, and 2) component-based architecture by mapping the traced classes into clusters. Execution traces are identified from the business tasks of end-users, i.e., the records of what classes have been executed during a user usage.

Static Analysis Component Identification from Applications. In [30], the authors used the Knowledge Discovery Meta-model to represent the software elements, at different levels of abstraction, and their structural dependencies. They used a vertical clustering algorithm to identify software components as groups of classes, and a horizontal clustering algorithm to recover the layered architecture view of the identified components. In [31], the authors proposed an approach to identify component interfaces of object-oriented components identified using reverse engineering approaches. Classes of a component are analyzed to structure the required and provided interfaces. They extract a set of methods invoked by other components. These methods are then grouped by means of formal concept analysis. Some other approaches relied on the analysis of several software applications at the same time to identify components cross these applications [32] [33] [16]. For example, in [16], components are identified as groups of classes

frequently presented in different applications and structural dependent.

Frequent Usage Patterns of API entities. Several approaches have been proposed to identify abstract reuse scenarios in terms of frequent usage patterns of API entities. These frequent usage patterns are not themselves direct reusable entities, but they help improving the reusability and understandability of APIs. We classified these existing approaches following their goals, usage-order consideration, the granularity of the API entities under the study, and the algorithm used to identify usage patterns. Generally the goal of usage-pattern identification approaches can be: (i) to provide examples to support recommendation systems [34] [35], (ii) to support the documentation of APIs at different levels of abstraction [34] [36] [5], (iii) to predict bugs resulted from incorrect usage scenarios [12], etc. The patterns were identified with respect to the order in which the API elements are used in some approaches [34] [36], while other approaches do not take into account such an order [12] [37] [5]. The granularity of the API elements that compose the identified patterns is at the method-level [34] [36] [5] or the class-level [1][37]. In [13], multi-level patterns are identified. As for the algorithms to identify the patterns, the approaches used association rules mining [37], frequent-pattern growth [1], clustering algorithms [36] [5] or a heuristic defined by the authors [34] [12]. Some approaches relied on a combination of many algorithms like Principle Component Analysis and clustering algorithm [35].

8 CONCLUSION AND FUTURE WORK

We propose an approach that aims to identify reusable software components based on the dynamic analysis of interactions between client applications and the targeted API. The approach relies on co-usage and objected-oriented dependencies to define relationships between classes of APIs. To do so, we execute usage scenarios of client applications of APIs to collect execution traces. Our approach packages groups of methods frequently appearing together in the collected traces as component provided interfaces, and their owner classes define the structure of components' implementation. These groups of methods are identified based on a graph-based clustering algorithm from the information in execution traces.

To evaluate our approach, we experimented with the DaCoPo benchmarks, focusing on three of its APIs and eight client applications. We considered pre-defined usage scenarios that are already provided by DaCoPo. The evaluation results show that the precision of our approach is up to 98%.

We will consider three future directions to investigate

- (1) We plan to transform the object-oriented implementation of the identified component provided interfaces into service-oriented interfaces to have truly reusable interfaces.
- (2) We want to extend the approach evaluation by considering more APIs and client applications to generalize the approach results.
- (3) We will also provide a visualization framework, that developers can use to identify reusable components of an API. They will have to load a set of client application using the API of interest, run the client applications while the framework is collecting the execution traces, and produces reusable components.

REFERENCES

- [1] Anas Shatnawi, Abdelhak-Djamel Seriai, Houari Sahraoui, and Zakarea Alshara. Reverse engineering reusable software components from object-oriented apis. *Journal of Systems and Software*, 131:442–460, 2017.
- [2] Minhaz F Zibran, Farjana Z Eishita, and Chanchal K Roy. Useful, but usable? factors affecting the usability of apis. In *2011 18th Working Conference on Reverse Engineering (WCRE)*, pages 151–155. IEEE, 2011.
- [3] Mohamed Aymen Saied, Hani Abdeen, Omar Benomar, and Houari Sahraoui. Could we infer unordered api usage patterns only using the library source code? In *Proceedings of the 2015 IEEE 23rd International Conference on Program Comprehension*, pages 71–81. IEEE Press, 2015.
- [4] Evan Moritz, Mario Linares-Vásquez, Denys Poshyvanyk, Mark Grechanik, Collin McMillan, and Malcom Gethers. Export: Detecting and visualizing api usages in large source code repositories. In *Proceedings of the 28th IEEE/ACM International Conference on Automated Software Engineering*, pages 646–651. IEEE Press, 2013.
- [5] Mohamed Aymen Saied and Houari Sahraoui. A cooperative approach for combining client-based and library-based api usage pattern mining. In *2016 IEEE 24th International Conference on Program Comprehension (ICPC)*, pages 1–10. IEEE, 2016.
- [6] Martin P Robillard and Robert Deline. A field study of api learning obstacles. *Empirical Software Engineering*, 16(6):703–732, 2011.
- [7] R Liguori and P Liguori. *Java 8 Pocket Guide*. "O'Reilly Media, Inc.", 2014.
- [8] Gias Uddin, Barthélémy Dagenais, and Martin P Robillard. Temporal analysis of api usage concepts. In *34th International Conference on Software Engineering*, pages 804–814. IEEE Press, 2012.
- [9] Mohamed Aymen Saied, Omar Benomar, Hani Abdeen, and Houari Sahraoui. Mining multi-level api usage patterns. In *2015 IEEE 22nd International Conference on Software Analysis, Evolution and Reengineering (SANER)*, pages 23–32. IEEE, 2015.
- [10] João Eduardo Montandon, Hudson Borges, Daniel Felix, and Marco Tulio Valente. Documenting apis with examples: Lessons learned with the apiminer platform. In *2013 20th Working Conference on Reverse Engineering (WCRE)*, pages 401–408. IEEE, 2013.
- [11] Mohamed Aymen Saied, Houari Sahraoui, and Bruno Dufour. An observational study on api usage constraints and their documentation. In *2015 IEEE 22nd International Conference on Software Analysis, Evolution and Reengineering (SANER)*, pages 33–42. IEEE, 2015.
- [12] Martin Monperrus, Marcel Bruch, and Mira Mezini. Detecting missing method calls in object-oriented software. In *European Conference on Object-Oriented Programming*, pages 2–25. Springer, 2010.
- [13] Hamzeh Eyal Salman. Identification multi-level frequent usage patterns from apis. *Journal of Systems and Software*, 130:42–56, 2017.
- [14] Oliviero Riganelli, Daniela Micucci, and Leonardo Mariani. Policy enforcement with proactive libraries. In *Proceedings of the 12th International Symposium on Software Engineering for Adaptive and Self-Managing Systems*, pages 182–192. IEEE Press, 2017.
- [15] Oliviero Riganelli, Daniela Micucci, Leonardo Mariani, and Yliès Falcone. Verifying policy enforcers. In *International Conference on Runtime Verification*, pages 241–258. Springer, 2017.
- [16] Anas Shatnawi and Abdelhak-Djamel Seriai. Mining reusable software components from object-oriented source code of a set of similar software. In *IEEE 14th International Conference on Information Reuse and Integration (IRI)*, pages 193–200. IEEE, 2013.
- [17] Seza Adjoayan, Abdelhak-Djamel Seriai, and Anas Shatnawi. Service identification based on quality metrics object-oriented legacy system migration towards soa. In *SEKE: Software Engineering and Knowledge Engineering*, pages 1–6. Knowledge Systems Institute Graduate School, 2014.
- [18] Abderrahmane Seriai, Salah Sadou, and Houari A Sahraoui. Enactment of components extracted from an object-oriented application. In *European Conference on Software Architecture*, pages 234–249. Springer, 2014.
- [19] Simon Allier, Salah Sadou, Houari Sahraoui, and Régis Fleurquin. From object-oriented applications to component-oriented applications via component-oriented architecture. In *9th Working IEEE/IFIP Conference on Software Architecture (WICSA)*, pages 214–223. IEEE, 2011.
- [20] Anas Shatnawi, Abdelhak Seriai, Houari Sahraoui, and Zakarea Al-Shara. Mining software components from object-oriented apis. In *International Conference on Software Reuse*, pages 330–347. Springer, 2015.
- [21] Aírel Pérez-Suárez, José F Martínez-Trinidad, Jesús A Carrasco-Ochoa, and José E Medina-Pagola. Oclustr: A new graph-based algorithm for overlapping clustering. *Neurocomputing*, 121:234–247, 2013.
- [22] Stephen M Blackburn, Robin Garner, Chris Hoffmann, Asjad M Khang, Kathryn S McKinley, Rotem Bentzur, Amer Diwan, Daniel Feinberg, Daniel Frampton, Samuel Z Guyer, et al. The dacapo benchmarks: Java benchmarking development and analysis. In *ACM Sigplan Notices*, volume 41, pages 169–190. ACM, 2006.
- [23] Zakarea Al-Shara, Abdelhak-Djamel Seriai, Chouki Tibermacine, Hinde Lilia Bouziane, Christophe Dony, and Anas Shatnawi. Materializing architecture recovered from oo source code in component-based languages. In *ECSCA: European Conference on Software Architecture*, 2016.
- [24] Zakarea Al-Shara, Abdelhak-Djamel Seriai, Chouki Tibermacine, Hinde Lilia Bouziane, Christophe Dony, and Anas Shatnawi. Migrating large object-oriented applications into component-based ones. In *GPCE: Generative Programming: Concepts and Experiences*, volume 51, pages 55–64, 2015.
- [25] OSGi Alliance. *Osgi service platform, release 3*. IOS Press, Inc., 2003.
- [26] Eric Bruneton, Thierry Coupaye, Matthieu Leclercq, Vivien Quéma, and Jean-Bernard Stefani. The fractal component model and its support in java. *Software: Practice and Experience*, 36(11-12):1257–1284, 2006.
- [27] Anas Shatnawi, Hamed Mili, Ghizlane El Boussaidi, Anis Boubaker, Yann-Gaël Guéhéneuc, Naouel Moha, Jean Privat, and Manel Abdellatif. Analyzing program dependencies in java ee applications. In *Proceedings of the 14th International Conference on Mining Software Repositories*, pages 64–74. IEEE Press, 2017.
- [28] Bas Cornelissen. Dynamic analysis techniques for the reconstruction of architectural views. In *14th Working Conference on Reverse Engineering (WCRE)*, pages 281–284. IEEE, 2007.
- [29] Philippe Dugerdil and David Sennhauser. Dynamic decision tree for legacy use-case recovery. In *28th Annual ACM Symposium on Applied Computing*, pages 1284–1291. ACM, 2013.
- [30] Ghizlane El Boussaidi, Alvine Boaye Belle, Stephane Vaucher, and Hamed Mili. Reconstructing architectural views from legacy systems. In *2012 19th Working Conference on Reverse Engineering (WCRE)*, pages 345–354. IEEE, 2012.
- [31] Abderrahmane Seriai, Salah Sadou, Houari Sahraoui, and Salma Hamza. Deriving component interfaces after a restructuring of a legacy system. In *2014 IEEE/IFIP Conference on Software Architecture (WICSA)*, pages 31–40. IEEE, 2014.
- [32] Anas Shatnawi, Abdelhak-Djamel Seriai, and Houari Sahraoui. Recovering software product line architecture of a family of object-oriented product variants. *Journal of Systems and Software*, 131:325–346, 2017.
- [33] Anas Shatnawi, Abdelhak Seriai, and Houari Sahraoui. Recovering architectural variability of a family of product variants. In *International Conference on Software Reuse*, pages 17–33. Springer, 2015.
- [34] J.E. Montandon, H. Borges, D. Felix, and M.T. Valente. Documenting apis with examples: Lessons learned with the apiminer platform. In *20th Working Conf. on Reverse Engineering (WCRE)*, pages 401–408, 2013.
- [35] G. Uddin, B. Dagenais, and M. P. Robillard. Temporal analysis of api usage concepts. In *Proc. of the 2012 Inter. Conf. on Software Engineering*, ICSE 2012, pages 804–814. Piscataway, NJ, USA, 2012. IEEE Press.
- [36] J. Wang, Y. Dang, H. Zhang, K. Chen, T. Xie, and D. Zhang. Mining succinct and high-coverage api usage patterns from source code. In *Proc. of the 10th Working Conf. on Mining Software Repositories*, MSR '13, pages 319–328. Piscataway, NJ, USA, 2013. IEEE Press.
- [37] M. Bruch, T. Schäfer, and M. Mezini. Fruit: Ide support for framework understanding. In *Proc. of the 2006 OOPSLA Workshop on Eclipse Technology eXchange*, eclipse '06, pages 55–59. New York, NY, USA, 2006. ACM.