# Recovering Architectural Variability of a Family of Product Variants

Anas Shatnawi<sup>1</sup>, Abdelhak Seriai<sup>1</sup>, and Houari Sahraoui<sup>2</sup>

<sup>1</sup> UMR CNRS 5506, LIRMM, University of Montpellier II, Montpellier, France shatnawi, seriai@lirmm.fr <sup>2</sup> DIRO, University of Montreal, Montreal, Canada sahraoui@iro.umontreal.ca

Abstract. A Software Product Line (SPL) aims at applying a preplanned systematic reuse of large-grained software artifacts to increase the software productivity and reduce the development cost. The idea of SPL is to analyze the business domain of a family of products to identify the common and the variable parts between the products. However, it is common for companies to develop, in an ad-hoc manner (e.g. clone and own), a set of products that share common functionalities and differ in terms of others. Thus, many recent research contributions are proposed to re-engineer existing product variants to a SPL. Nevertheless, these contributions are mostly focused on managing the variability at the requirement level. Very few contributions address the variability at the architectural level despite its major importance. Starting from this observation, we propose, in this paper, an approach to reverse engineer the architecture of a set of product variants. Our goal is to identify the variability and dependencies among architectural-element variants at the architectural level. Our work relies on Formal Concept Analysis (FCA) to analyze the variability. To validate the proposed approach, we experimented on two families of open-source product variants; Mobile Media and Health Watcher. The results show that our approach is able to identify the architectural variability and the dependencies.

**Keywords:** Product line architecture · architecture variability · architecture recovery · product variants · reverse engineering · source code · object-oriented.

# 1 Introduction

A Software Product Line (SPL) aims at applying a pre-planned systematic reuse of large-grained software artifacts (e.g. components) to increase the software productivity and reduce the development cost [1–3]. The main idea behind SPL is to analyze the business domain of a family of products in order to identify the common and the variable parts between these products [1, 2]. In SPL, the variability is realized at different levels of abstraction (e.g. requirement and design). At the requirement level, it is originated starting from the differences in users' wishes, and does not carry any technical sense [2] (e.g. the user needs *camera* and *WIFI* features in the phone). At the design level, the variability starts to have more details related to technical senses to form the product architectures. These technical senses are described via Software Product Line Architecture (SPLA). Such technical senses are related to which components compose the product (e.g. video recorder, and photo capture components), how these components interact through their interfaces (e.g. video recorder provides a video stream interface to media store), and what topology forms the architectural configuration (i.e. how components are composited and linked) [2].

Developing a SPL from scratch is a highly costly task since this means the development of the domain software artifacts [1]. In addition, it is common for companies to develop a set of software product variants that share common functionalities and differ in terms of other ones. These products are usually developed in an ad-hoc manner (e.g. clone and own) by adding or/and removing some functionalities to an existing software product to meet the requirement of a new need [4]. Nevertheless, when the number of product variants grows, managing the reuse and maintenance processes becomes a severe problem [4]. As a consequence, it is necessary to identify and manage variability between product variants as a SPL. The goal is to reduce the cost of SPL development by first starting it from existing products and then being able to manage the reuse and maintenance tasks in product variants using a SPL. Thus, many research contributions have been proposed to re-engineer existing product variants into a SPL [5,6]. Nevertheless, existing works are mostly focused on recovering the variability in terms of features defined at the requirement level. Despite the major importance of the SPLA, there is only two works aiming at recovering the variability at the architectural level [7, 8]. These approaches are not fullyautomated and rely on the domain knowledge which is not always available. Also, they do not identify dependencies among the architectural elements. To address this limitation, we propose in this paper an approach to automatically recover the architecture of a set of software product variants by capturing the variability at the architectural level and the dependencies between the architectural elements. We rely on Formal Concept Analysis (FCA) to analyze the variability. In order to validate the proposed approach, we experimented on two families of open-source product variants; Mobile Media and Health Watcher. The evaluation shows that our approach is able to identify the architectural variability and the dependencies as well.

The rest of this paper is organized as follows. Section 2 presents the background needed to understand our proposal. Then, in Section 3, we present the recovery process of SPLA. Section 4 presents the identification of architecture variability. Then, Section 5 presents the identification of dependencies among architectural-element variants. Experimental evaluation of our approach is discussed in section 6. Then, the related work is discussed in Section 7. Finally, concluding remarks and future directions are presented in section 8.

## 2 Background

# 2.1 Component-Based Architecture Recovery from Single Software: ROMANTIC Approach

In our previous work [9, 10],  $ROMANTIC^3$  approach has been proposed to automatically recover a component-based architecture from the source code of

<sup>&</sup>lt;sup>3</sup> *ROMANTIC*: Re-engineering of Object-oriented systeMs by Architecture extractioN and migraTIon to Component based ones.

an existing object-oriented software. Components are obtained by partitioning classes of the software. Thus each class is assigned to a unique subset forming a component. ROMANTIC is based on two main models. The first concerns the object-to-component mapping model which allows to link object-oriented concepts (e.g. package, class) to component-based ones (e.g. component, interface). A component consists of two parts; internal and external structures. The internal structure is implemented by a set of classes that have direct links only to classes that belong to the component itself. The external structure is implemented by the set of classes that have direct links to other components' classes. Classes that form the external structure of a component define the component interface. Fig. 1 shows the object-to-component mapping model. The second main model proposed in *ROMANTIC* is used to evaluate the quality of recovered architectures and their architectural-element. For example, the quality-model of recovered components is based on three characteristics; composability, autonomy and specificity. These refer respectively to the ability of the component to be composed without any modification, to the possibility to reuse the component in an autonomous way, and to the fact that the component implements a limited number of closed functionalities. Based on these models, ROMANTIC defines a fitness function applied in a hierarchical clustering algorithm [9, 10] as well as in search-based algorithms [11] to partition the object-oriented classes into groups, where each group represents a component. In this paper, ROMANTIC is used to recover the architecture of a single object oriented software product.



Fig. 1. Object-to-component mapping model

#### 2.2 Formal Concept Analysis

Formal Concept Analysis (FCA) is a mathematical data analysis technique developed based on lattice theory [12]. It allows the analysis of the relationships between a set of objects described by a set of attributes. In this context, maximal groups of objects sharing the same attributes are called formal concepts. These are extracted and then hierarchically organized into a graph called a concept lattice. Each formal concept consists of two parts. The first allows the representation of the objects covered by the concepts called the extent of the concept. The second allows the representation of the set of attributes shared by the objects belonging to the extent. This is called the intent of the concept. Concepts can be linked through sub-concept and super-concept relationship [12] where the lattice defines a partially ordered structure. A concept A is a sub-concept of the super-concept B, if the extent of the concept B includes the extent of the concept B.





Table 1. Formal context

Fig. 2. Lattice of formal context in Table 1

The input of FCA is called a formal context. A formal context is defined as a triple K = (O, A, R) where O refers to a set of objects, A refers to a set of attributes and R is a binary relation between objects and attributes. This binary relation indicates to a set of attributes that are held by each object (i.e.  $R \subseteq OXA$ ). Table 1 shows an example of a formal context for a set of bodies of waters and their attributes. An X refers to that an object holds an attribute.

As stated before, a formal concept consists of extent E and intent I, with E a subset of objects O ( $E \subseteq O$ ) and I a subset of attributes A ( $I \subseteq A$ ). A pair of extent and intent (E, I) is considered a formal concept, if and only, if E consists of only objects that shared all attributes in I and I consists of only attributes that are shared by all objects in E. The pair ("river, lake", "inland, natural, constant") is an example of a formal concept of the formal context in Table 1. Fig. 2 shows the concept lattice of the formal context presented in Table 1.

# 3 Process of Recovering Architectural Variability

The goal of our approach is at recovering the architectural variability of a set of product variants by statically analyzing their object-oriented source code. This is obtained by identifying variability among architectures respectively recovered from each single product. We rely on *ROMANTIC* approach to extract the architecture of a single product. This constitutes the first step of the recovery process. Architecture variability is related to architectural-elements variability, i.e. component, connector and configuration variability. In our approach, we focus only on component and configuration variability<sup>4</sup>. Fig. 3 shows an example of architecture variability based on component and configuration variability. In this example, there are three product variants, where each one diverges in the set of component constituting its architecture as well as the links between the components. *CD Reader* and *CD Reader / Writer* represent variants of one component. We identify component variants based on the identification of components providing similar functionalities. This is the role of the second step of

<sup>&</sup>lt;sup>4</sup> Most of architectural description languages do not consider connector as a first class concept.

the recovery process. Configuration variability is represented in terms of presence/absence of components on the one hand (e.g. *Purchase Reminder*), and presence/absence of component-to-component links on the other hand (e.g. the link between *MP3 Decoder / Encoder* and *CD Reader / Writer*). We identify configuration variability based on both the identification of core (e.g. *Sound Source*) and optional components (e.g. *Purchase Reminder*) and links between these components. In addition, we capture the dependencies and constraints among components. This includes, for example, require constraints between optional components. We rely on FCA to identify these dependencies. These are mined in the fourth step of the recovery process. Fig. 4 shows these steps.

# 4 Identifying the Architecture Variability

The architecture variability is mainly materialized either through the existence of variants of the same architectural element (i.e. component variants) or through the configuration variability. In this section, we show how component variants and configuration variability are identified.

## 4.1 Identifying Component Variants

The selection of a component to be used in an architecture is based on its provided and required services. The provided services define the role of the component. However, other components may provide the same, or similar, core services. Each may also provide other specific services in addition to the core ones. Considering these components, either as completely different or as the same, does not allow the variability related to components to be captured. Thus, we consider them as component variants. We define component variants as a set of components providing the same core services and differ concerning few secondary ones. In Fig. 3, *MP3 Decoder* and *MP3 Decoder* / *Encoder* are component variants.



Fig. 3. An example of architecture variability

Fig. 4. The process of architectural variability recovery

We identify component variants based on their similarity. Similar components are those sharing the majority of their classes and differing in relation to some others. Components are identified as similar based on the strength of similarity links between their implementing classes. For this purpose, we use cosine similarity metric [13] where each component is considered as a text document composed of the names of its classes. We use a hierarchical clustering algorithm [13] to gather similar components into clusters. It starts by considering components as initial leaf nodes in a binary tree. Next, the two most similar nodes are grouped into a new one that forms their parent. This grouping process is repeated until all nodes are grouped into a binary tree. All nodes in this tree are considered as candidates to be selected as groups of similar components. To identify the best nodes, we use a depth first search algorithm. Starting from the tree root to find the cut-off points, we compare the similarity of the current node with its children. If the current node has a similarity value exceeding the average similarity value of its children, then the cut-off point is in the current node. Otherwise, the algorithm continues through its children. The results of this algorithm are clusters where each one is composed of a set of similar components that represent variants of one component.

#### 4.2 Identifying Configuration Variants

The architectural configuration is defined based on the list of components composing the architecture, as well as the topology of the links existing between these components. Thus the configuration variability is related to these two aspects; the lists of core (mandatory) and optional components and the list of core and optional links between the selected components.

**Identification of component variability:** To identify mandatory and optional components, we use Formal Concept Analysis (FCA) to analyze architecture configurations. We present each software architecture as an object and each member component as an attribute in the formal context. In the concept lattice, common attributes are grouped into the root while the variable ones are hierarchically distributed among the non-root concepts.



Fig. 5. A lattice example of similar configurations

Fig. 5 shows an example of a lattice for three similar architecture configurations. The common components (the core ones) are grouped together at the root concept of the lattice (the top). In Fig. 5 *Com1* and *Com4* are the core components present in the three architectures. By contrast, optional components are represented in all lattice concepts except the root. e.g., according to the lattice of Fig. 5, *Com2* and *Com5* present in *Arch1* and *Arch2* but not in *Arch3*.

Identification of component-link variability: A component-link is defined as a connection between two components where each connection is the abstraction of a group of method invocation, access attribute or inheritance links between classes composing these components. In the context of configuration variability, a component may be linked with different sets of components. A component may have links with a set of components in one product, and it may have other links with a different set of components in another product. Thus the component-link variability is related to the component variability. This means that the identification of the link variability is based on the identified component variability. For instance, the existence of a link A-B is related to the selection of a component A and a component B in the architecture. Thus considering a core link (mandatory link) is based on the occurrence of the linked components, but not on the occurrence in the architecture of products. According to that, a core link is defined as a link occurring in the architecture configuration as well the linked components are selected. To identify the component-link variability, we proceed as follows. For each architectural component, we collect the set of components that are connected to it in each product. The intersection of the sets extracted from all the products determines all core links for the given component. The other links are optional ones.

## 5 Identifying Architecture Dependencies

The identification of component and component-link variability is not enough to define a valid architectural configuration. This also depends on the set of dependencies (i.e. constraints) that may exist between all the elements of the architecture. For instance, components providing antagonism functionalities have an exclude relationship. Furthermore, a component may need other components to perform its services. Dependencies can be of five kinds: alternative, OR, AND, require, and exclude dependencies. To identify them we rely on the same concept lattice generated in the previous section.

In the lattice, each node groups a set of components representing the intent and a set of architectural configurations representing the extent. The configurations are represented by paths starting from their concepts to the lattice concept root. The idea is that each object is generated starting from its node up going to the top. This is based on sub-concept to super-concept relationships (c.f. Section 2.2). This process generates a path for each object. A path contains an ordered list of nodes based on their hierarchical distribution; i.e. sub-concept to super-concept relationships). According to that, we propose extracting the dependencies between each pair of nodes as follows:

- Required dependency. This constraint refers to the obligation selection of a component to select another one; i.e. component B is required to select component A. Based on the generated lattice, we analyze all its nodes by

identifying parent-to-child relation (i.e. top to down). Thus node A requires node B if node B appears before node A in the lattice, i.e., node A is a subconcept of the super-concept corresponding to node B. In other words, to reach node A in the lattice, it is necessary to traverse node B. For example, if we consider lattice of the Fig. 5, Com6 requires Com2 and Com5 since Com2 and Com5 are traversed before Com6 in all paths including Com6and linking root node to object nodes.

- Exclude and alternative dependencies. Exclude dependency refers to the antagonistic relationship; i.e. components A and B cannot occur in the same architecture. This relation is extracted by checking all paths linking root to all leaf nodes in the lattice. A node is excluded with respect to another node if they never appear together in any of the existing paths; i.e. there is no sub-concept to super-concept relationship between them. This means that there exists no object exists containing both nodes. For example, if we consider lattice of Fig. 5, *Com6* and *Com7* are exclusives since they never appear together in any of the lattice paths.

Alternative dependency generalizes the exclude one by exclusively selecting only one component from a set of components. It can be identified based on the exclude dependencies. Indeed, a set of nodes in the lattice having each an exclude constraint with all other nodes forms an alternative situation.

- AND dependency. This is the bidirectional version of the REQUIRE constraint; i.e. component A requires component B and vice versa. More generally, the selection of one component among a set of components requires the selection of all the other components. According to the built lattice, this relation is found when a group of components is grouped in the same concept node in the lattice; i.e. the whole node should be selected and not only a part of its components. For example if we consider lattice of the Fig. 5, Com2 and Com5 are concerned with an AND dependency.
- OR dependency. When some components are concerned by an OR dependency, this means that at least one of them should be selected; i.e. the configuration may contain any combination of the components. Thus, in the case of absence of other constraints any pair of components is concerned by an OR dependency. Thus pairs concerned by required, exclude, alternative, or AND dependencies are ignored as well as those concerned by transitive require constraints; e.g. *Com6* and *Com7* are ignored since they are exclusives. Algorithm 1 shows the procedure of identifying groups of OR dependency.

# 6 Experimentation and Results

Our experimentation aims at showing how the proposed approach is applied to identify the architectural variability and validating the obtained results. To this end, we applied it on two case studies. We select two sets of product variants. These sets are Mobile Media<sup>5</sup> (MM) and Health Watcher<sup>6</sup> (HW). We select these products because they were used in many research papers aiming at addressing

 $<sup>^5</sup>$  Available at: http://ptolemy.cs.iastate.edu/design-study/#mobilemedia.

<sup>&</sup>lt;sup>6</sup> Available at: http://ptolemy.cs.iastate.edu/design-study/#healthwatcher.

Algorithm 1: Identifying OR-Groups

the problem of migrating product variants into a SPL. Our study considers 8 variants of MM and 10 variants of HW. MM variants manipulate music, video and photo on mobile phones. They are developed starting from the core implementation of MM. Then, the other features are added incrementally for each variant. HW variants are web-based applications that aim at managing health records and customer complaints. The size of each variant of MM and HW, in terms of classes, is shown in Table 2. We utilize *ROMANTIC* approach [9] to extract architectural components from each variant independently. Then, the components derived from all variants are the input of the clustering algorithm to identify component variants. Next, we identify the architecture configurations of the products. These are used as a formal context to extract a concept lattice. Then, we extract the core (mandatory) and optional components as well as the dependencies among optional-component.

In order to evaluate the resulted architecture variability, we study the following research questions:

- RQ1: Are the identified dependencies correct? This research question goals at measuring the correctness of the identified component dependencies.
- RQ2: What is the precision of the recovered architectural variability? This research question focuses on measuring the precision of the resulting architecture variability. This is done by comparing it with a preexisted architecture variability model.

Table 2. Size of MM variants and HW ones

Name	1	2	3	4	5	6	7	8	9	10	Avg.
MM	25	34	36	36	41	50	60	64	Х	Х	43.25
HW	115	120	132	134	136	140	144	148	160	167	136.9

#### 6.1 Results

Table 3 shows the results of component extraction from each variant independently, in terms of the number of components, for each variant of MM and HW.



Fig. 6. The concept lattice of HW architecture configurations

The results show that classes related to the same functionality are grouped into the same component. The difference in the numbers of the identified components in each variant has resulted from the fact that each variant has a different set of user's requirements. On average, a variant contains 6.25 and 7.7 main functionalities respectively for MM and HW.

 Table 3. Comp. extraction results

 Table 4. Comp. variants identification

Name	1	2	3	4	5	6	7	8	9	10	Avg.	Total
MM	3	5	5	5	7	$\overline{7}$	9	9	Х	Х	6.25	50
HW	6	7	9	10	7	9	8	8	7	6	7.7	77

Table 4 summarizes the results of component variants in terms of the number of components having variants (NOCV), the average number of variants of a component (ANVC), the maximum number of component variants (MXCV) and the minimum number of component variants (MNCS). The results show that there are many sets of components sharing the most of their classes. Each set of components mostly provides the same functionality. Thus, they represent variants of the same architectural component. Table 5 presents an instance of 6 component variants identified from HW, where X means that the corresponding class is a member in the variant. By analyzing these variants, it is clear that these components represent the same architectural component. In addition to that, we noticed that there are some component variants having the same set of classes in multiple product variants.

The architecture configurations are identified based on the above results. Table 6 shows the configuration of MM variants, where X means that the component is a part of the product variants. The results show that the products are similar in their architectural configurations and differ considering other ones. The reason behind the similarity and the difference is the fact that these products are common in some of their user's requirements and variable in some others. These architecture configurations are used as a formal context to extract the concept

Class Name	Variant 1	Variant 2	Variant 3	Variant 4	Variant 5	Variant 6
BufferedReader	Х	Х	Х	Х	Х	Х
ComplaintRepositoryArray	Х	Х	Х	Х	Х	Х
ConcreteIterator	Х	Х	Х	Х	Х	Х
DiseaseRecord	Х					
<b>IIteratorRMITargetAdapter</b>	Х	Х	Х	Х	Х	Х
IteratorRMITargetAdapter	Х	Х	Х	Х	Х	Х
DiseaseType		Х				
InputStreamReader	Х	Х	Х	Х	Х	Х
Employee		Х		Х		
InvalidDateException			Х	Х	Х	Х
IteratorDsk	Х	Х	Х	Х	Х	Х
PrintWriter	Х	Х	Х		Х	Х
ObjectNotValidException			Х		Х	Х
RemoteException	Х	Х		Х		
PrintStream			Х		Х	Х
RepositoryException	Х	Х				
Statement	Х	Х	Х	Х	Х	Х
Throwable	Х	Х		Х		
HWServlet					Х	
Connection					Х	Х

Table 5. Instance of 6 component variants

lattice. We use the Concept Explorer<sup>7</sup> tool to generate the concept lattice. Due to limited space, we only give the concept lattice of HW (c.f. Fig. 6). In Table 7, the numbers of core (mandatory) and optional components are given for MM and HW. The results show that there are some components that represent the core architecture, while some others represent delta (optional) components.

**Table 6.** Arch. configuration for allMM variants

Variant No	o. Com1	Com2	Com3	Com4	Com5	Com6	Com7	Com8	Com9	Com1(	Com1:	Com1:	Com1:	$Com1_{2}$
1	X		Х		Х						_			-
2	X		Х		Х			Х	Х					
3	X		Х		Х			Х	Х					
4			Х		Х			Х	Х		Х			
5	X		Х		Х			Х	Х		Х	Х		
6			Х		Х		Х	Х	Х		Х	Х		
7		Х		Х	Х		Х	Х	Х	Х	Х	Х		
8		Х		Х	Х	Х	Х	Х		Х			Х	Х

**Table 7.** Mandatory and optionalcomponents

Product Name	MM	HW
Mandatory	1	2
Optional	13	16

The results of the identification of optional-component dependencies are given in Table 8 (*Com 5* is excluded since it is a mandatory component). For conciseness, the detailed dependencies among components are only shown for MM only. The dependencies are represented between all pairs of components in MM (where R= Require, E= Exclude, O= OR, RB= Required By, TR= Transitive Require, TRB = Transitive Require By, and A = AND). Table 9 shows a summary of MM and HW dependencies between all pairs of components. This includes the number of direct require constrains (NRC), the number of exclude ones (NE), the number of AND groups (NOA), and the number of OR groups (NO). Alternative constrains is represented as exclude ones. The results show that there are dependencies among components that help the architect to avoid

<sup>&</sup>lt;sup>7</sup> Presentation of the Concept Explorer tool is available in [14].

creating invalid configuration. For instance, a design decision of AND components indicates that these components depend on each other, thus, they should be selected all together.

Table 8. Component dependencies

	C1	C2	C3	C4	C6	C7	C8	C9	C10	C11	C12	C13	C14
Com1	X		R		Е	Е				0	Е	Е	E
Com2		Х	Е	Α	RB	R	TR		Α			RB	RB
Com3	RB	E	Х	Е	Е		0		Е			Е	Е
Com4		A	Е	Х	RB	R	TR		Α			RB	RB
Com6	Е	R	Е	R	Х	TR	TR	Е	R	Е	Е	Α	Α
Com7	E	RB		RB	TRB	Х	R	0	RB			TRB	TRB
Com8		TRB	0	TRB	TRB	RB	Х	RB	TRB	TRB	TRB	TRB	TRB
Com9					Е	0	R	Х		RB	TRB	Е	E
Com 10		Α	Е	Α	RB	R	TR		Х			RB	RB
Com11	0				Е		TR	R		Х	RB	Е	E
Com12	E				Е		TR	TR		R	Х	Е	Е
Com 13	E	R	Е	R	Α	TR	TR	Е	R	Е	Е	Х	A
Com14	E	R	Е	R	Α	TR	TR	Е	R	Е	Е	Α	Х

**Table 9.** Summarization ofMM and HW dependencies

Name	NDR	NE	NA	NO
MM	17	20	6	3
HW	18	62	3	11

To the best our knowledge, there is no architecture description language supporting all kinds of the identified variability. The existing languages are mainly focused on modeling component variants, links and interfaces, while they do not support dependencies among components such as AND-group, OR-group, and require. Thus, on the first hand, we use some notation presented in [15] to represent the concept of component variants and links variability. On the other hand, we propose some notation inspired from feature modeling languages to model the dependencies among components. For the purpose of understandability, we document the resulting components by assigning a name based on the most frequent tokens in their classes' names. Figure 7 shows the architectural variability model identified for MM variants, where the large boxes denote to design decisions (constraints). For instance, core architecture refers to components that should be selected to create any concrete product architecture. In MM, there is one core components manipulating the base controller of the product. This component has two variants. A group of Multi Media Stream, Video Screen Controller, and Multi Screen Music components represents an AND design decision.

**RQ1:** Are the identified dependencies correct? The identification of component dependencies is based on the occurrence of components. e.g., if two components never selected to be included in a concrete product architecture, we consider that they hold an exclude relation. However, this method could provide correct or incorrect dependencies. To evaluate the accuracy of this method, we manually validate the identified dependencies. This is based on the functionalities provided by the components. For instance, we check if the component functionality requires the functionality of the required component and so on. The results show that 79% of the required dependencies are correct. As an example of a correct relation is that *SMS Controller* requires *Invalid Exception* as it performs an input/output operations. On the other hand, it seems that *Image Util* does not require *Image Album Vector Stream*. Also, 63% of the exclude constrains are correct. For AND and OR dependencies, we find that 88% of AND groups are correct, while 42% of OR groups are correct. Thus, the precision of identifying dependencies is 68% in average.



Fig. 7. Architectural variability model for MM

RQ2: What is the precision of the recovered architectural variability? In our case studies, MM is the only case study that has an available architecture model containing some variability information. In [16], the authors presented the aspect oriented architecture for MM variants. This contains information about which products had added components, as well as in which product a component implementation was changed (i.e. component variants). We manually compare both models to validate the resulting model. Fig. 8 shows the comparison results in terms of the total number of components in the architecture model (TNOC), the number of components having variants (NCHV), the number of mapped components in the other model (NC), the number of unmapped components in the other model (NUMC), the number of optional components (NOC) and the number of mandatory ones (NOM). The results show that there are some variation between the results of our approach and the pre-existed model. The reason behind this variation is the idea of compositional components. For instance, our approach identifies only one core component compared to 4 core components in the other model. Our approach grouped all classes related to the controller components together in one core components. On the other hand, the other model divided the controller component into Abstract Controller, Album Data, Media



*Controller*, and *Photo View Controller* components. In addition, the component related to handling exceptions is not mentioned in the pre-existed model at all.

Fig. 8. The results of the MM validation

## 7 Related Work

In this section, we discuss the contributions that have been proposed in two research directions; recovering the software architecture of a set of product variants and variability management.

In [7], an approach aiming at recovering SPLA was presented. It identifies component variants based on the detection of cloned code among the products. However, the limitation of this approach is that it is a semi-automated, while our approach is fully automated. Also, it does not identify dependencies among the components. In [8], the authors presented an approach to reconstruct Home Service Robots (HSR) products into a SPL. Although this approach identifies some architectural variability, but it has some limitation compared to our approach. For instance, it is specialized on the domain of HSR as the authors classified, at earlier step, the architectural units based on three categories related to HSR. These categories guide the identification process. In addition, the use of feature modeling language (hierarchical trees) to realize the identified variability is not efficient as it is not able to represent the configuration of architectures. Domain knowledge plays the main role to identify the architecture of each single product and the dependencies among components. In some cases, domain knowledge is not always available. The authors in [6] proposed an approach to reverse engineering architectural feature model. This is based on the software architect's knowledge, the architecture dependencies, and the feature model that is extracted based on a reverse engineering approach presented in [5]. The idea, in [6], is to take the software architect's variability point of view in the extracted feature model (i.e. still at the requirement level); this is why it is named architecture feature model. However, the major limitations of this approach are firstly that the software architect is not available in most cases of legacy software, and secondly that the architecture dependencies are generally missing as well. In [5], the authors proposed an approach to extract the feature model. The input of the extraction process is feature names, feature descriptions and dependencies among features. Based on this information, they recover ontological constraints (e.g. feature groups) and cross tree constrains. A strong assumption behind this approach is that feature names, feature descriptions, and dependencies among features are available. In [17], the authors use FCA to generate a feature model. The input of their approach is a set of feature configurations. However, the extraction of the feature model elements is based on NP-hard problems (e.g. set cover to identify or groups). Furthermore, architecture variability is not taken into account in this approach. In [18], an approach was presented to visually analyze the distribution of variability and commonality among the source code of product variants. The analysis includes multi-level of abstractions (e.g. line of code, method, class, etc.). This aims to facilitate the interpretation of variability distribution, to support identifying reusable entities. In [19], the authors presented an approach to extract reusable software components from a set of similar software products. This is based on identifying similarity between components identified independently from each software. This approach can be related only to the first step of our approach.

## 8 Conclusion

In SPLA, the variability is mainly represented in terms of components and configurations. In the case of migrating product variants to a SPL, identifying the architecture variability among the product variants is necessary to facilitate the software architect's tasks. Thus, in this paper, we proposed an approach to recover the architecture variability of a set of product variants. The recovered variability includes mandatory and optional components, the dependencies among components (e.g. require, etc.), the variability of component-links, and component variants. We rely on FCA to analyze the variability. Then, we propose two heuristics. The former is to identify the architecture variability. The latter is to identify the architecture dependencies. The proposed approach is validated through two sets of product variants derived from Mobile Media and Health Watcher. The results show that our approach is able to identify the architectural variability and the dependencies as well.

There are three aspects to be considered regarding the hypothesis of our approach. Firstly, we identify component variants based on the similarity between the name of classes composing the components, i.e., classes that have the same name should have the same implementation. While in some situations, components may have very similar set of classes, but they are completely unrelated. Secondly, dependencies among components are identified based on component occurrences in the product architectures. Thus, the identified dependencies maybe correct or incorrect. Finally, the input of our approach is the components independently identified form each product variants using ROMAN-TIC approach. Thus, the accuracy of the obtained variability depends on the accuracy of ROMANTIC approach.

Our future research will focus on migrating product variants into component based software product line, the mapping between the requirements' variability (i.e. features) and the architectures' variability, and mapping between components' variability and component-links' variability.

#### References

- 1. Clements, P., Northrop, L.: Software product lines: practices and patterns. Addison-Wesley Reading (2002)
- Pohl, K. and Böckle, G. and Van Der Linden, F.: Software product line engineering. Springer Berlin Heidelberg (2005)
- Tan, L. and Lin, Y. and Ye, H.: Quality-oriented software product line architecture design. Journal of Software Engineering & Applications. 5(7), 472–476 (2012)
- Rubin, J., Chechik, M.: Locating distinguishing features using diff sets. In: IEEE/ACM 27th Inter. Conf. on ASE, pp. 242–245. (2012)
- She, S., Lotufo, R., Berger, T., Wasowski, A., Czarnecki, K.: Reverse engineering feature models. In: Proc. of 33rd ICSE, pp. 461–470. (2011)
- Acher, M., Cleve, A., Collet, P., Merle, P., Duchien, L., Lahire, P.: Reverse engineering architectural feature models. In: Software Architecture. LNCS, vol. 6903, pp. 220–235. Springer, Heidelberg (2011)
- Koschke, R., Frenzel, P., Breu, A.P., Angstmann, K.: Extending the reflexion method for consolidating software variants into product lines. Software Quality Journal. 17(4), 331–366 (2009)
- Kang, K.C., Kim, M., Lee, J., Kim, B.: Feature-oriented re-engineering of legacy systems into product line assets - a case study. In: Software Product Lines. LNCS, vol. 3714, pp. 45–56. Springer, Heidelberg (2005)
- Kebir, S., Seriai, A.D., Chardigny, S., Chaoui, A.: Quality-centric approach for software component identification from object-oriented code. In: Proc. of WICSA/ECSA, pp. 181–190. (2012)
- Chardigny, S., Seriai, A., Oussalah, M., Tamzalit, D.: Extraction of componentbased architecture from object-oriented systems. In: Proc. of 7th WICSA, pp. 285– 288. (2008)
- Chardigny, S., Seriai, A.D., Oussalah, M., Tamzalit, D.: Search-based extraction of component-based architecture from object-oriented systems. In: 2nd ECSA. LNCS, vol. 5292, pp. 322–325. Springer, Heidelberg (2008)
- 12. Ganter, B., Wille, R.: Formal concept analysis. WISSENSCHAFTLICHE ZEITSCHRIFT-TECHNISCHEN UNIVERSITAT DRESDEN. 47, 8–13 (1996)
- Han, J., Kamber, M., Pei, J.: Data mining: concepts and techniques. Morgan kaufmann (2006)
- Yevtushenko, A.S.: System of data analysis "concept explorer". (In Russian) Proc. of the 7th National Conf. on Artificial Intelligence (KII), vol. 79, pp. 127–134. (2000)
- 15. Hendrickson, S.A., van der Hoek, A.: Modeling product line architectures through change sets and relationships. In: Proc. of the 29th ICSE, pp. 189–198. (2007)
- Figueiredo, E. and Cacho, N. and Sant'Anna, C. and Monteiro, M. and Kulesza, U. and Garcia, A. and Soares, S. and Ferrari, F. and Khan, S., et al.: Evolving software product lines with aspects. In: Proc. of 30th ICSE, pp. 261–270. (2008)
- Ryssel, U., Ploennigs, J., Kabitzsch, K.: Extraction of feature models from formal contexts. In: Proc. of 15th SPLC, pp. 1–4. (2011)
- Duszynski, S., Knodel, J., Becker, M.: Analyzing the source code of multiple software variants for reuse potential. In: Proc. of WCRE, pp. 303–307. (2011)
- Shatnawi, A., Seriai, A.D.: Mining reusable software components from objectoriented source code of a set of similar software. In: IEEE 14th Inter. Conf. on Information Reuse and Integration (IRI), pp. 193–200. (2013)