

# Conceptual Code Mining

## Mining for Source-Code Regularities with Formal Concept Analysis

Kim Mens\*

*Département d'Ingénierie Informatique (INGI)  
Université catholique de Louvain (UCL)  
Place Sainte Barbe 2, B-1348 Louvain-la-Neuve, Belgium*

Tom Tourwé

*Centrum voor Wiskunde en Informatica (CWI)  
P.O. Box 94079, NL-1090 GB Amsterdam, The Netherlands*

---

### Abstract

Understanding the conceptual structure of large software systems, whether it is for software understanding or reengineering purposes, is a nontrivial task. In particular, knowing *where to start* the comprehension process is more difficult than it seems, especially when a system is large and complex and time is scarce. We propose an approach to mine a system's source code automatically and efficiently for relevant concepts of interest, which we refer to as *source-code regularities*: what concerns are addressed in the code, what patterns, programming idioms and conventions have been adopted, and where and how they are implemented. We use formal concept analysis to do the actual source-code mining, and then filter, classify and combine the results to present them in a format that is more convenient to a software engineer. We applied a prototype tool that implements this approach to several small to medium-sized Smalltalk applications. For each of these, the tool discovered several interesting source-code regularities. Although the tool and approach can still be improved in many ways, the tool does already provide useful results when having a *first contact* with a system. The obtained results also illustrate the relevance and feasibility of using formal concept analysis as a technique for source code mining.

*Key words:* Source-code mining, formal concept analysis, software classification.

---

\* Corresponding author.

*Email addresses:* Kim.Mens@info.ucl.ac.be (Kim Mens), Tom.Tourwe@cwi.nl (Tom Tourwé).

*URLs:* <http://www.info.ucl.ac.be/people/cvmens.html> (Kim Mens),

## 1 Introduction

“Where do I start?” is one of the most crucial questions to be answered when embarking on an important reengineering task. Demeyer et al. [7] even devote an entire chapter of their book on object-oriented reengineering to this “first contact” with some unfamiliar code. Apart from being confronted with a system that is completely new, for this first contact phase typically only a few days are allocated, during which an initial and reasonably accurate assessment of the system needs to be made.

Although a variety of software reengineering support tools exist, few tools are dedicated to this first contact phase, when it is not really clear yet what exactly we are looking for. Most tools need a minimum of information and a specific goal before they can be applied. Hence, in this first contact phase a software reengineer often performs better by relying on informal techniques such as “chat with maintainers”, “read all code in one hour”, “skim documentation”, “interview during demo” or “do a mock installation” [7].

Nevertheless, to complement and enhance these techniques, we would like to have some more automated support for this phase. Therefore, we investigate whether the technique of *formal concept analysis* [11], with known applications in data analysis and knowledge processing, is applicable at this early stage. The essence of our contribution lies not in the idea of applying FCA to source code, but in our particular choice of elements and properties for the FCA algorithm and how we filtered and classified the discovered concepts, in order to mine a system’s source code for relevant *source-code regularities*, in a way that is independent of the actual system being analyzed.

Although the proposed approach can still be improved in many ways, and in spite of its apparent simplicity, it allows us to mine Smalltalk source code for many interesting regularities, like design patterns, programming idioms and conventions. It also allows us to discover certain opportunities for refactoring, as well as some features of which the implementation is spread throughout the source code. Most of the discovered information provides a good starting point for understanding the source code in more detail.

The remainder of this paper is structured as follows. Section 2 briefly introduces the mathematical technique of formal concept analysis. In Section 3 we explain our approach and how we used formal concept analysis to mine the source code for relevant concepts of interest. Section 4 briefly presents the tool and Section 5 gives an overview of the experiments we conducted and the source-code regularities we discovered. Sections 6 and 7 discuss related and future work. We conclude the paper in Section 8.

---

<http://www.cwi.nl/~tourwe> (Tom Tourwé).

## 2 Formal Concept Analysis

Formal concept analysis (FCA) [11] is a branch of lattice theory that can be used to identify meaningful groupings of *elements* that have common *properties*.<sup>1</sup> The FCA algorithm takes as input a relation, or Boolean table,  $T$  between a (potentially large, but finite) set of elements and a set of properties of those elements. An example of such a table is given in Table 1, in which different programming languages and properties are related. A cross in a table cell means that the programming language in the corresponding row has the property of the corresponding column.

Table 1

Programming languages and their supported programming paradigms.

Progr. language	OO	Functional	Logic	Static typing	Dynamic typing
Java	X	-	-	X	-
Smalltalk	X	-	-	-	X
C++	X	-	-	X	-
Scheme	-	X	-	-	X
Prolog	-	-	X	-	X

Taking such a table  $T$  as input, the FCA algorithm determines *maximal* groups of elements and properties, called *concepts*, such that :

- each element of the group shares the properties,
- every property of the group holds for all of its elements,
- no other element outside the group has those same properties, and
- no other property outside the group holds for all elements in the group.

Intuitively, a *concept* corresponds to a maximal ‘rectangle’ in the table  $T$ , where we allow arbitrary permutations of the table’s rows and columns.

All concepts are ordered into a *concept lattice*, an example of which is depicted in Figure 1. The lattice’s bottom concept contains those elements that have all properties. Since there is no such programming language in our example, that concept contains no elements. Similarly, the top concept contains those properties that hold for all elements. Again, there is no such property. Other concepts represent related groups of programming languages, such as the concept  $(\{ Java, C++ \}, \{ static\ typing, object\ oriented \})$ , which groups all statically-typed object-oriented languages.

---

<sup>1</sup> As in Arévalo et al. [2], in this paper we prefer to use the terms *element* and *property* instead of *object* and *attribute* used in traditional FCA literature, because these latter terms already have a very specific meaning in OO software development.

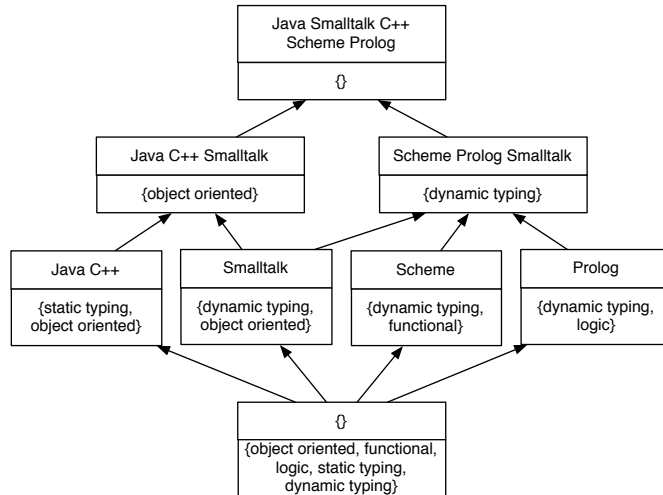


Fig. 1. Concept lattice corresponding to Table 1.

For more details on formal concept analysis we refer to [11]. The next section explains the details of our approach to use FCA for source code mining.

### 3 Mining for Source-Code Regularities

When applying FCA for mining Smalltalk source code, we first have to choose the elements and properties to compute the concept lattice (§3.1). When computing the lattice (§3.2), lots of concepts are produced, many of which are irrelevant or redundant. Therefore, we filter the discovered concepts (§3.3) and classify (§3.4), combine and annotate them (§3.5) in a way that is more relevant to a software engineer.

The novelty of our contribution is not in the idea of applying FCA to source code. What is more important is our particular choice of elements, properties, filters and analyzers, and how these allow us to discover interesting source-code regularities in a way that is independent of the considered application.<sup>2</sup>

#### 3.1 Generate FCA elements and properties

Since our goal in this paper is to mine Smalltalk source code, as FCA *elements* we chose source-code entities like classes, methods and method parameters. Other source-code entities such as variables, method protocols, class categories, bundles, packages and namespaces were not considered in this initial

<sup>2</sup> and, so we believe, even largely independent of the considered programming language, even though we did not yet validate that claim with a real experiment.

experiment, in order to avoid cluttering the results.

As *properties* we take simple substrings of the names of these source-code entities. As such, the discovered concepts will group entities with similar names. Our motivation for choosing these properties, is that in Smalltalk in particular and in many other programming languages, programmers often rely on naming conventions to reveal their intentions and to implement certain programming idioms and design patterns. Keeping the properties simple has the additional benefit that they can be generated and manipulated efficiently.

Nevertheless, to limit the number of properties, we do not consider all possible substrings. Instead, we split class, method and parameter names in substrings according to the capitals and other separators occurring in them. In addition, we discard substrings with little conceptual meaning or that are used too often, such as ‘with’, ‘from’, ‘the’, ‘object’, as well as substrings that are too small (i.e., less than 3 characters). We also ignore colons, plurals and the difference in case when comparing substrings. For example, the properties associated with a class `QuotedCodeConstant` are the substrings ‘quoted’, ‘code’ and ‘constant’. The properties corresponding to a method named `#unifyWithDelayedVariable:inEnv:myIndex:hisIndex:inSource:` are ‘unify’, ‘delayed’, ‘variable’, ‘env’, ‘index’ and ‘source’.

### 3.2 Compute the concept lattice

Applying an FCA algorithm to the elements and properties generated in the previous step, results in a large *concept lattice* of several hundreds to thousands of concepts (cf. Table 2, Section 5, which shows some quantitative results of applying our approach on five different Smalltalk applications.)

For now, let it suffice to give a single illustrative example of a kind of concepts that is discovered by the FCA algorithm: *accessing methods*. Indeed, since both accessor and mutator methods are named after an instance variable, they share a same substring. E.g., the following concept we discovered in the *SOUL* application groups the `#callStack` accessor and `#callStack:` mutator methods of the `callStack` instance variable defined in the `Environment` class:

```
Environment >> callStack
  ^ callStack
Environment >> callStack: aStack
  callStack := aStack
```

They are grouped based on the properties ‘call’ and ‘stack’ that are shared by these methods, and by no other methods in *SOUL*.

### 3.3 Filter the concepts

As we will see in Section 5 (cf. third column *#raw* of Table 2), the number of concepts discovered by FCA, before applying any filtering, is of the same order of magnitude as the number of considered elements. This would imply that a software engineer needs to look at a significant number of concepts in order to try and understand the source code. Luckily, however, there is a lot of redundancy and noise in the discovered concepts. To reduce some of this noise, we apply some simple *filters*.

A first filter ignores all concepts that contain two or less elements, since these concepts are generally too small to provide relevant information. Note that this filter discards most *accessing method* concepts, since these typically contain only two elements: an accessor and a mutator method. However, since accessing methods are rather fine-grained, since there are a lot of them, and since they can be inspected with standard browsers easily or they can be retrieved with more dedicated tools, we don't mind that they get discarded.

A second filter ignores all concepts that share only one property (substring). Although this filter may discard some interesting concepts, it does throw away many more irrelevant concepts. We think that, in a 'first contact' setting, getting a quick and focussed idea of certain source-code regularities is more important than getting a precise list of *all* possible regularities. (A nice improvement of this filter would be to discard those concepts of which the properties 'cover' only a small fraction, based on some threshold, of the name of the elements. As such, the filter becomes relative to the size of the elements' names.)

Whereas these two generic filters are independent of the kinds of elements being analyzed, our third filter is more targeted. It discards concepts that contain only classes (with a similar name) in the same hierarchy. These concepts typically do not provide very useful information (except if we want to discover exactly which naming convention these classes are relying upon), since classes belonging to the same hierarchy often have similar names.

### 3.4 Classify the filtered concepts

Being mere sets of elements (classes and methods) and properties (substrings of the elements' names), the concepts that remain after filtering are rather unstructured. Therefore, we reorganize the concepts automatically in a way that is more easy for a software engineer to analyze and interpret.

More precisely, we flatten the concept lattice and *classify* the concepts in a way that makes more sense to the software engineer. We visualize the classifications

of concepts directly in the StarBrowser [16]. We distinguish 3 main groups of concepts (in fact, the classification we used was much finer-grained than the one presented here, but space restrictions prohibit us to explain it entirely):

- (1) *Single class concepts* group concepts of which all elements are methods (or parameters of those methods) belonging to a single class;
- (2) *Hierarchy concepts* have a larger scope as they group classes, methods and parameters of those methods, that belong to a single class hierarchy;
- (3) For *crosscutting concepts* we explicitly require that at least two different class hierarchies are involved. We do this by verifying that the most specific common superclass of the considered classes is `object` and that none of the methods in the concept are defined on the `object` class itself (which would be a degenerate case of a *hierarchy concept*).

Such a taxonomy helps a software engineer in several ways. Depending on what he/she is trying to understand or looking for — a single class, a hierarchy or a crosscutting concern — he/she will prefer one of the above classifications over the other. In addition, knowing that a certain concept belongs to a given classification helps him/her to better understand that concept. For example, knowing that a concept containing several methods with exactly the same name belongs to the *hierarchy concepts* classification allows him/her to qualify those methods as *polymorphic methods*.

### 3.5 Combine and annotate concepts

By organizing the concepts in a taxonomy like the above, the structure of the lattice is lost. Concepts that were nearby in the lattice (e.g., that were in a subconcept relationship) will not necessarily belong to the same classification, and vice versa. As a consequence, since there is a lot of overlap between concepts that are nearby in the lattice, when reorganizing the concepts this may lead to redundancy among concepts that get classified into *different* classifications. Whenever possible, however, when nearby concepts in the concept lattice are put in the *same* classification, we automatically reconstruct part of the original structure of the lattice, in order to reduce redundancy. More specifically, we recombine highly overlapping concepts into a single nested one.

In addition, we automatically regroup and annotate the concepts belonging to each classification, in order to present them in a way that is more convenient to the software engineer: different concepts related to the same class(es) are combined, methods are annotated with the classes they belong to, and concepts are annotated with their properties.

## 4 DelfSTof, our Conceptual Code Mining tool

We developed a prototype tool that implements the approach outlined above, and presents the discovered concepts in a way that is easy to use and manipulate. The tool is called *DelfSTof*, coining the Dutch word ‘delfstof’, which designates the result of a delving process. (In English, the first meaning of the verb “to delve” is “to make careful investigation for facts, knowledge, etc.”. Coincidentally, the pronunciation of the word “delfstof” sounds like the English “delve stuff” which is indeed what the tool does.) We capitalize the letters “ST” because the tool is implemented completely in Smalltalk and analyzes Smalltalk source code. Nevertheless, the approach and tool can be generalized easily to analyze the source code of other languages as well; we intend to generalize it for mining Java source code in the near future.

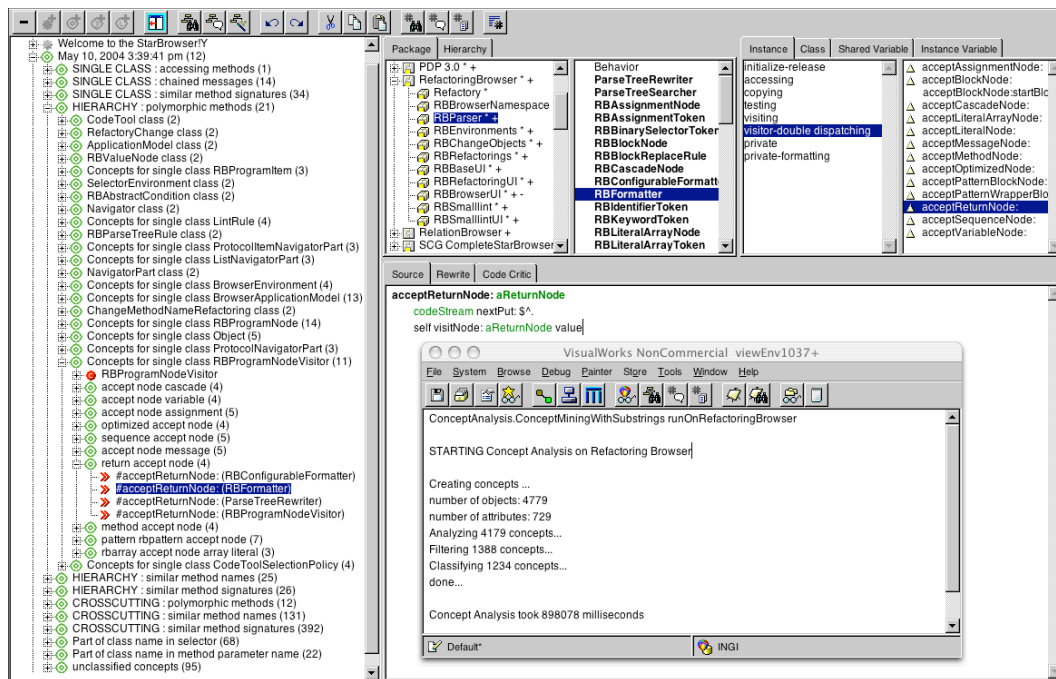


Fig. 2. DelfSTof : Discovered concepts for the Refactoring Browser.

Without going into all details, the tool consists of an efficient FCA algorithm, a set of filters, and a set of ‘analyzers’ that are in charge of the classification, combination and annotation of concepts. The resulting taxonomy of concepts is visualized with the *StarBrowser* [16]. A screenshot of the tool, which is essentially a *StarBrowser* plugin, is presented in Figure 2.

## 5 Discovered Source-Code Regularities

We applied our tool to five different cases, as summarized by Table 2 below.

Table 2

Quantitative results of FCA applied to some Smalltalk applications.

Case	#elements	#properties	#raw	#filtered	time (sec)
DelfSTof	756 (135)	237	617	126	5
StarBrowser	731 (52)	352	740	115	7
SOUL	1469 (111)	434	1188	281	22
CodeCrawler	1370 (93)	477	1419	327	24
Refactoring Browser	4779 (271)	729	4179	1234	414

Every row corresponds to a different application of which we mined the source code for regularities. *SOUL* is an interpreter for a Prolog-like language and *DelfSTof* is our own conceptual code mining tool. We chose these two applications because we know their implementation well, which allowed us to better assess the relevance of the discovered results. Both the *StarBrowser* and the *Refactoring Browser* are advanced Smalltalk browsers and *CodeCrawler* is a language-independent reverse engineering tool which combines metrics and software visualization. These cases were chosen because we wanted to better understand their implementation anyway, in the context of another experiment and tool we are currently working on.

The column *#elements* gives an indication of the size of each considered case as it equals the total number of classes plus methods in that case. (The number between brackets is the number of classes.) Note that we did not store method parameters as separate elements since they are part of the method signature.

The column *#properties* shows the number of substrings that have been generated from the considered elements. We observe that the number of properties is always a factor 2 to 4, in the case of the Refactoring Browser even almost a factor 7, *less* than the number of elements. This is a good sign as it implies that a significant amount of the properties are actually shared by the elements.

The third column *#raw* shows the raw number of concepts discovered by FCA. Column *#filtered* shows how many concepts remain after having applied the simple filters that were explained in §3.3. We observe that, after filtering, there remain about 4 to 7 times less concepts than the number of considered elements. We still think that this is a bit too much, especially for larger cases, but we will come back to this discussion in Section 7.

For all considered cases, the time of computation — which includes all steps explained in Section 3 and not only the computation of the concept lattice — was acceptable (ranging from a few seconds to a few minutes), although it increases in a non-linear way with the number of considered objects.

The remainder of this section discusses some of the regularities we discov-

ered when manually analyzing the results of our FCA experiments. We could refine the taxonomy of §3.4 to classify explicitly and automatically some of these regularities as well. However, a certain trade-off needs to be made, since this requires more automated analysis up front and thus slows down the tool. Also, we want to keep the tool sufficiently general so that it still can be applied to other languages (maybe even non OO languages). Nevertheless, we already extended our tool so that it classifies automatically *accessing methods*, *polymorphic methods*, *chained messages*, and *delegating methods*(cf. §5.1).

### 5.1 Programming idioms

**Accessing methods**, as explained before, are easily recognized by our tool because of the naming convention they rely on. However, because there typically exist a lot of accessing methods, and because they do not really offer very valuable information, we discarded most of them by applying a filter that requires at least 3 source-code entities per concept.

**Polymorphic methods** are another source-code regularity that is readily recognized by our tool, since polymorphic methods have exactly the same name. Consequently, they are grouped together in a concept. In addition, if there are several polymorphic methods for a same class hierarchy, these will all be grouped together automatically in a single combined concept.

For example, in the source code of the *Refactoring Browser* we discovered a concept containing 4 methods named `#acceptReturnNode`, implemented on different subclasses of `RBProgramNodeVisitor`. This is a typical example of polymorphism of which many examples can be found in any OO application. In fact, for the class `RBProgramNodeVisitor` alone many more polymorphic method concepts were discovered, as can be seen from Figure 2.

*Polymorphic methods* across class hierarchies are equally interesting to detect as they may trigger interesting refactorings or tell us something about possible multiple inheritance problems the original developers encountered.

**Chained messages** are concepts that group a method together with some of its auxiliary methods in the same class. These chains are recognized by FCA since the auxiliary methods often have a name that is similar to that of the originating method, though sometimes slightly longer and taking an extra parameter. E.g., in the *CodeCrawler* application the class `CCMetricsChooserDialog` implements a method `#applyChosenMetrics`, which calls an auxiliary method `#detectChosenMetrics`, which in turns calls `#assignChosenMetricsTo:`. These 3 methods share the substrings ‘chosen’ and ‘metrics’.

**Delegating methods** delegate responsibility by calling a method with the same name. Our tool discovered some interesting sets of delegating methods with a similar name that all belonged to the same class. The presence of many such delegating methods in a single class may indicate that the class is a *Decorator* [10].

## 5.2 Code duplication

By closely inspecting the discovered concepts, we also detected several cases of copy and paste reuse: several concepts contain methods that not only have a similar name, but a similar implementation as well. This may seem logical, since methods that implement similar behavior can be expected to have similar names. However, from an implementation point of view, this duplicated code could just as well be factored out and reused.

For example, in *CodeCrawler*'s `CEVModelHistory` class we discovered a concept with 2 methods `#predecessorModelNameOfModel:` and `#predecessorModelNameOfModelNamed:` which had nearly the same implementation. Since one of them is no longer being used, we assume that the original developer(s) created one of the methods by copying it from the other one, then replaced all calls to the old one to the new one, but in the end forgot to remove the old method.

In general, particular reasons for duplication may become clear by looking at the classification of the concepts:

- a developer who was not aware that a method implementing the desired behavior was already defined, may accidentally implement a method with a similar name and behavior. Such methods are often grouped in concepts classified as *hierarchy concepts*, because the method that already implements the desired behavior is probably not implemented by the class the developer is looking at, but by one of its sub- or superclasses.
- a method was needed whose behavior differed slightly from the behavior already implemented by an existing method, and this behavior was copied and adapted slightly, without extracting the common code into a separate method (or merely deleting the original version if it is no longer needed, such as in the example of the `CEVModelHistory` class above). Concepts containing such duplicated methods tend to be classified as *single class concepts*, because such duplication typically occurs inside a single class.
- the duplicated behavior could not be factored out into a single piece of code, and thus could not be reused. This is mostly due to the fact that the duplicated code occurs in classes defined in different class hierarchies. As such, these concepts are classified as *crosscutting concepts*.

### 5.3 Design patterns

As many design patterns [10] use certain naming conventions, it is no surprise that they are detected by our tool. For example, the *Visitor* pattern uses the convention that each *visit* method defined by a *visitor* class encodes the name of the class being visited. Since they clearly share some substrings, our tool will group a class and its corresponding visit methods inside a single concept.

*SOUL* uses a *Visitor* design pattern in order to perform a variety of operations on logic terms, such as copying and renaming them. The terms are represented as subclasses of the `AbstractTerm` hierarchy, while the visitor hierarchy is defined by the `SimpleTermVisitor` class hierarchy. Our tool recovered this design pattern instance in two concepts:

- (1) A combined concept in classification *hierarchy concepts* which groups all *polymorphic methods* in the `SimpleTermVisitor` hierarchy, that implement the behavior to be executed when a particular term is visited. The concept consists of a several sub-concepts, each of which contains all methods defined by subclasses of class `SimpleTermVisitor`, dealing with one particular term. For example, one such concept is defined by the properties ‘visit’ and ‘compound’, and contains various implementations of the `compoundVisit:` method, defined in the `SimpleTermVisitor` hierarchy and responsible for implementing behavior associated to a `CompoundTerm` object. More specifically, the concept consists of four `compoundVisit:` methods, implemented in the classes `SimpleTermVisitor`, `CompoundTermRenamingVisitor`, `CopyingVisitor` and `LexicalAddressVisitor`.
- (2) The second concept is also a *hierarchy concept* and contains all `accept` methods defined by subclasses of `AbstractTerm`. These methods are responsible for calling the appropriate method, corresponding to the term being visited, in the supplied visitor object. They are grouped based on the ‘visitor’ and ‘accept’ substring properties. This is an example of a concept that takes into account both the method’s name and the name of its formal parameter, since the `accept:` methods always define a formal parameter named `aVisitor`.

Note that the *polymorphic methods* depicted in Figure 2 are also a part of a visitor pattern, used in the *Refactoring Browser* implementation.

As several other design patterns use similar naming conventions, we also detected occurrences of the *Abstract Factory*, *Builder*, *Observer* and *Decorator* design patterns. Due to space limitations we cannot give examples of these.

#### 5.4 Relevant domain concepts

Frequently occurring properties give a good idea of what the important concepts in the application or problem domain are. This information is very useful to understand the domain, and to provide a common vocabulary which can be used to talk with maintainers. For example when applying our FCA tool to the source code of *DelfSTof* itself, we found several concepts with properties like ‘concept’, ‘attribute’, ‘analyze(r)’, ‘filter’ or ‘classification’, which are indeed important concepts in the domain of formal concept analysis.

#### 5.5 Opportunities for refactoring

Besides revealing interesting source-code regularities, our approach can identify opportunities for refactorings [9] that improve the source code quality.

An obvious opportunity for refactoring is to get rid of some of the code duplication that was detected (§5.2). The way a concept containing duplicated methods is classified, can provide useful hints about which refactorings to apply. If the concept is classified as a *single class concept*, the duplication occurs in a single class, and an *Extract method* refactoring is appropriate. If the concept occurs in the *hierarchy concepts* classification, a combination of *Extract method* and *Pull up method* refactorings is probably more appropriate. Of course, if the duplication is caused by having copied a method that is no longer used, it suffices to simply remove that method.

Concepts that were recognized as *polymorphic methods* can also be inspected for refactoring opportunities, to ensure that polymorphism is well-implemented:

- We could check whether all classes in a class hierarchy understand the polymorphic method. If not, we may need to add an additional one.
- A polymorphic method might also be implemented by several subclasses of a particular superclass, but not by that superclass itself. In that case, a *Pull up method* or *Add class* refactoring may be appropriate to define the methods in the superclass, or to insert an intermediate superclass.

A particular example we discovered in *SOUL* is the `#updateRepositories` method, which is only defined separately in subclasses `RepositoryBrowser`, `SoulQueryBrowser` and `SoulClauseBrowser` but not in their common superclass `ApplicationModel`. Introducing an intermediate superclass here might be appropriate.

In the *CodeCrawler* case, we also found several examples where the *same* polymorphic method appeared in several subclasses, but not in their common superclass. However, in most of these examples the code in one sub-hierarchy did not seem to be used, which made us believe that the code was moved from one part in the hierarchy to another, but that the developer(s)

forgot to remove the original code.

- *Crosscutting polymorphic methods* are also suspicious. For example, we discovered that in *SOUL*, the `#usesPredicate:multiplicity:` method is implemented in both the `AbstractTerm` and `HornClause` hierarchies. Smalltalk, which has dynamic typing, still allows classes of these hierarchies to be used polymorphically. A static type system, as in Java, would prohibit such polymorphic use. In any case, this situation may call for a refactoring.

Concepts that represent design pattern instances can also be scanned for particular design flaws. For example, the *Visitor* design pattern requires that each visitor class defines an appropriate `visit` method and, vice versa, that each element class defines an `accept` method that calls the appropriate `visit` method. The concepts that identify instances of the *Visitor* design pattern can be used to inspect the implementation in a quick and straightforward way, and verify whether these constraints are adhered to.

## 6 Related Work

The use of FCA in software engineering is not new. Snelting et al. [17] use FCA to reengineer C++ class hierarchies, while Arévalo et al. [2] analyze object-oriented framework reuse using FCA. Closer to our work are the techniques by Tonella et al. [19] and Ducasse et al. [6]. The former use FCA to detect instances of design patterns in source code. Since they specifically tune the FCA algorithm for detecting such instances, they are not able to detect other kinds of regularities, as our approach does. The latter use FCA to reveal the structure of single classes only. They partition the methods of a class, according to the fields these methods use, and then use the concept lattice to visualize and understand the structure of that class. Tilley et al. [18] provide an overview of the use of FCA for several other software engineering purposes.

A large number of tools to verify the quality of the source code of an application exists. The spectrum ranges from very simple tools that detect basic coding errors [15], over specialized clone detection tools [3,4,8,13], to tools that detect high-level bad smells [12,20,21] and propose appropriate refactorings. Other tools exist that are capable of detecting high-level structures in source code, such as coding conventions and design patterns [1,5,14,22].

The main difference between these tools and ours, is that our approach requires no a priori knowledge. Most of these existing tools, however, rely on the fact that design pattern implementations follow *particular* naming conventions and guidelines. Our approach is not targeted to detecting a specific kind of regularities, but is able to detect a variety of regularities that reveal bad code (bad smells, duplication), good code (design patterns, programming

idioms), and opportunities for refactoring. This is particularly useful in a first contact setting, whereas in a later phase, when the code is better understood, a more directed tool is preferable.

## 7 Future work

An important topic of future work is to further *improve the filtering* of the concepts discovered by FCA, so as to reduce the remaining redundancy in the discovered concepts. The problem is that this redundancy occurs between concepts that are classified in different categories. As was briefly mentioned in §3.3, this redundancy is a consequence of having flattened the concept lattice. By doing so we lost some important dependencies between concepts. But exactly this information may be useful to get rid of the redundancy. Therefore, we propose to keep the lattice as an internal representation, so that advanced filters can take advantage of it to resolve the remaining redundancy.

Since the time of computation of our tool increases in a non-linear way with the number of considered objects, this may pose problems regarding the scalability of the approach. However, we do not think that it is a good idea to apply the approach to *very* large amounts of source code, since the tool assumes that certain naming convention are adhered to in a consistent way. This is unlikely for very large cases. In such a context it would be better to apply the approach multiple times on several smaller pieces of which we know they are more or less independent and have been developed by a same development team. As a side-effect, this will also resolve the problem with the time of computation.

## 8 Conclusion

In this paper we proposed a fairly efficient tool, capable of mining an application's source code for meaningful and interesting regularities. The tool combines formal concept analysis with filtering and classification techniques, in order to provide simple and effective views on important parts of the source code. To validate the approach we applied the tool on a number of small to medium-sized Smalltalk applications. In spite of relying on nothing more than similarity of names of source-code entities, we discovered regularities such as programming idioms, code duplication, design patterns and domain concepts, as well as interesting opportunities for refactoring. Although the approach can still be improved in many ways, in particular to further reduce the redundancy, we do believe it can be very useful in a *first contact* context, in which little a priori knowledge about the source code is available and time is scarce.

## 9 Acknowledgements

We thank T. Mens, R. Wuyts, R. Riquelme, S. González and G. Arévalo as well as the anonymous reviewers of the ECOOP'2004 *workshop on object-oriented reengineering* and the AOSD'2004 workshop on *foundations of aspect languages* for their feedback on earlier versions of this paper. We thank F. Spiessens for his efficient implementation of the FCA algorithm in Smalltalk.

## References

- [1] H. Albin-Amiot, P. Cointe, Y.-G. Guéhéneuc, and N. Jussien. Instantiating and Detecting Design Patterns: Putting Bits and Pieces Together. In *Proc. Int. Conf. on Automated Software Engineering*, 2001.
- [2] G. Arévalo and T. Mens. Analysing object-oriented application frameworks using concept analysis. *LNCS*, 2426:53–63, September 2002.
- [3] B. S. Baker. On Finding Duplication and Near-Duplication in Large Software Systems. In *Proc. Second IEEE Working Conference on Reverse Engineering*, pages 86–95, July 1995. Received IEEE Outstanding Paper Award.
- [4] I. D. Baxter, A. Yahin, L. Moura, M. Sant'Anna, and L. Bier. Clone detection using abstract syntax trees. In *Proc. Int. Conf. on Software Maintenance (ICSM'98)*, pages 368–377. IEEE Computer Society Press, 1998.
- [5] K. Brown. Design reverse-engineering and automated design pattern detection in smalltalk. Master's thesis, University of Illinois at Urbana Champaign, 1996.
- [6] U. Dekel and Y. Gil. Revealing Class Structure with Concept Lattices. In *Proc. 10th Working Conference on Reverse Engineering*, 2003.
- [7] S. Demeyer, S. Ducasse, and O. Nierstrasz. *Object-Oriented Reengineering Patterns*. Morgan Kaufmann and DPunkt, 2002.
- [8] S. Ducasse, M. Rieger, and S. Demeyer. A language independent approach for detecting duplicated code. In H. Yang and L. White, editors, *Proc. Int. Conf. Software Maintenance*, pages 109–118. IEEE Computer Society Press, September 1999.
- [9] M. Fowler. *Refactoring: Improving the Design of Existing Code*. Addison-Wesley, 1999.
- [10] E. Gamma, R. Helm, R. Johnson, and J. Vlissides. *Design Patterns: Elements of Reusable Object-Oriented Software*. Addison Wesley, Massachusetts, 1994.
- [11] B. Ganter and R. Wille. *Formal Concept Analysis: Mathematical Foundations*. Springer-Verlag, 1999.

- [12] Y. Kataoka, M. D. Ernst, W. G. Griswold, and D. Notkin. Automated Support for Program Refactoring using Invariants. In *Proc. Int. Conf. on Software Maintenance*, pages 736–743, 2001.
- [13] R. Komondoor and S. Horwitz. Using Slicing to Identify Duplication in Source Code. In *Proc. of the 8th International Symposium on Static Analysis*. Springer-Verlag, 2001.
- [14] K. Mens, I. Michiels, and R. Wuyts. Supporting Software Development through Declaratively Codified Programming Patterns. *Journal on Expert Systems with Applications*, December 2002.
- [15] S. Paul and A. Prakash. A Framework for Source Code Search using Program Patterns. *IEEE Transactions on Software Engineering*, 20(6), June 1994.
- [16] W. Roel and D. Stéphane. Unanticipated integration of development tools using the classification model. *Journal of Computer Languages, Systems and Structures*, 30:pp. 63–77, 2003.
- [17] G. Snelling and F. Tip. Reengineering Class Hierarchies Using Concept Analysis. In *Proc. ACM SIGSOFT Symposium on the Foundations of Software Engineering*, 1998.
- [18] T. Tilley, R. Cole, P. Becker, and P. Eklund. A Survey of Formal Concept Analysis Support for Software Engineering Activities. In *Proc. 1st International Conference on Formal Concept Analysis*, 2003.
- [19] P. Tonella and G. Antoniol. Inference of object oriented design patterns. *Journal of Software Maintenance - Research and Practice*, 13(5):309 – 330, September - October 2001.
- [20] T. Tourwé and T. Mens. Identifying Refactoring Opportunities Using Logic Meta Programming. In *Proc. 7th European Conf. on Software Maintenance and Reengineering*, pages 91 – 100, Benvento, Italy, 2003. IEEE Computer Society.
- [21] E. van Emden and L. Moonen. Java quality assurance by detecting code smells. In *Proc. 9th Working Conference on Reverse Engineering*. IEEE Computer Society Press, October 2002.
- [22] R. Wuyts. Declarative Reasoning about the Structure of Object-Oriented Systems. In *Proc. TOOLS USA '98, IEEE Computer Society Press*, pages 112–124, 1998.