

Effective API Navigation and Reuse

Awny Alnusair, Tian Zhao
Department of Computer Science
University of Wisconsin-Milwaukee, USA
{alnusair, tzhao}@uwm.edu

Eric Bodden
Department of Computer Science
Technische Universität Darmstadt, Germany
bodden@acm.org

Abstract

Programmers can frequently benefit from reusing existing program code. However, most reuse libraries come with few source code examples that demonstrate how the library at hand should be used. We have developed a source-code recommendation approach for constructing and delivering relevant code snippets that programmers can use to complete a certain programming task. Our approach is semantic-based, relying on an explicit ontological representation of source code. We argue that such representation opens new doors for an improved recommendation mechanism that ensures relevancy and accuracy. Many existing recommendation systems require an existing repository of relevant code samples. However, for many libraries, such a repository does not exist. In our approach, we instead utilize points-to analysis to infer precise type information of library components. We have backed our approach with a tool that has been tested on multiple software frameworks. The obtained results are promising and demonstrate the effectiveness of our approach.

1. Introduction

Code reuse has many benefits when developing, enhancing, and evolving large-scale software systems. Programmers can reuse code in many ways. One common way is to access libraries of reusable components, or to plug into application frameworks. Unfortunately, many libraries and frameworks are not intuitive to use. While some exceptional pieces of software may be documented well, it is often the case that libraries and frameworks lack informative API documentation, and lack sufficient and effective source code examples that would explain a particular library feature. When such an example does exist, it can be very helpful: programmers can often simply copy the example into the current project and then adapt it to the new context, thus enabling a particular API programming task to be completed rather quickly.

In this paper, we present an approach for automatic source-code recommendation. Our approach is based on the idea that many programming tasks require programmers to compose a chain of method calls that convert a given source object of some particular type to a target object of some other type. Thus, the programmer needs to answer object-instantiation queries of the form (*Source object* \Rightarrow *Destination object*). In the special case where the *Source object* is not specified, the object instantiation-problem is reduced to either a simple constructor invocation or a static method invocation. For illustration, consider a programmer trying to reuse Jena¹, an open-source Java application framework for building Semantic Web applications. At some point, the programmer wishes to obtain an Ontology Class view of a Resource (e.g., ontological property) identified by a given URI. This view can be used, for example, to determine if an Ontology Class with the same URI exists in the Ontology Model. A programming task of this kind can be seen as an object-instantiation task of the form (*Property* \Rightarrow *OntClass*). The following code snippet shows a sample solution for this query:

```
Property pr = ..... ;
Model model = pr.getModel();
OntModel ontModel = (OntModel) model;
OntClass ontClass =
    ontModel.getOntClass(pr.getURI());
```

In an ideal world, even for a programmer who is new to Jena, accomplishing this conceptually simple task should not be time consuming. In practice however, identifying the proper sequence of method calls, as well as the required type cast within this chain, require the programmer to have a substantial knowledge of the framework's structure. Even worse, some complex call sequences require getting a handle of an object by invoking a static method from yet a totally different class (cf. snippet in Section 3.1). Furthermore, many application frameworks do not provide type-specific methods that may free the developer from using

¹<http://jena.sourceforge.net/>

downward type casts. This complicates the process of composing code snippets, as casts must be inserted to make the snippets compile. Our proposal to automatic code recommendation tackles these issues effectively by providing an ontology-based representation of source code structures that captures type information about possible concrete types.

An ontology is an explicit specification of a conceptualization [4]. That is, an ontology provides means to formally and explicitly describe concepts, objects, properties and other entities in a domain of discourse, and to describe the relationships that hold among these concepts. We thus use ontology formalisms to represent software assets by building a knowledge base that is automatically populated with instances representing source code artifacts. Our approach uses this knowledge base to effectively identify and retrieve relevant code snippets.

Besides addressing knowledge representation issues, our approach improves on the existing state-of-the-art in the following ways:

- ◊ Unlike other recommendation approaches, our approach neither requires a repository of sample code to mine for snippets, nor do we require our tool to be backed by a source-code search engine to obtain these samples. These requirements have been identified as one of the major limitations of current recommendation systems [10]. Instead, we solely rely on analyzing the framework’s or library’s API.
- ◊ Since the structure of an API usually contains too little information to obtain useful code snippets that require special features (e.g. determining the legality of a type cast), we use interprocedural points-to analysis to enrich our knowledge base with information about the possible runtime behavior of the API.
- ◊ We provide a context-sensitive approach that analyzes the user’s code when constructing, ranking, and delivering potential code-snippet candidates. For instance, the sample solution discussed above may not be the best candidate if there was an object of type `OntModel` already visible in the current context.
- ◊ Similar to other approaches, we traverse a graph representation of source code to find a path from the source to the destination object. However, since our graph is based on an ontology model, it is enriched with additional data that guides the search and helps us obtain more precise results.

2 Related work

Researchers have proposed many code-mining techniques, all of which tackle the recommendation problem

from different perspectives. Data Mining-based techniques try to reveal usage patterns of program components from a corpus of existing code examples. This is usually accomplished by extracting association rules which incorporate taxonomies of inheritance relationships [8], or by applying frequent-sequence mining and clustering techniques [16] to extract API methods that are frequently invoked in sequence. Such data mining techniques often suffer from scalability and rule-complexity issues. Traditional information-retrieval techniques, on the other hand, circumvent some of these complexity issues. They allow users to formulate keyword queries to retrieve source code samples ranked based on the match between the query and the obtained name-based indices [11] or latent semantic-based indices [9]. Due to the nature of the schemes and the keyword-based search employed, traditional keyword-based recommendation systems are usually very imprecise.

Some other approaches do in fact recommend personalized code snippets when queried. These approaches base their recommendation on analyzing a large corpus of sample client code collected using Google Code Search (GCS) (PARSEWeb [13]), or by searching in a pre-populated local repository (Strathcona [5], Prospector [7], and XSnippet [12]). Strathcona for example, is a recommendation tool that uses heuristics to match the structure of the code under development (structural context) to the structure of the code in a source code repository. PARSEWeb, Prospector, and XSnippet, on the other hand, are more focused on answering specific object instantiation queries.

Although these approaches take important steps in the right direction, we believe that there are fundamental issues related to the mechanisms used for data gathering, data processing, and most importantly, data representation. Firstly, with the exception of Prospector, all approaches that we mentioned rely on a hard-to-find repository populated with client code that expresses good usages of the framework. Prospector on the other hand, does analyze API signatures for the most part, but still relies on a repository to handle special features such as downcasts. For tools that collect code samples from the web, these samples are usually incomplete fragments that cannot be analyzed precisely. Secondly, traditional knowledge-representation mechanisms and hard-coded heuristics affect the quality of the retrieved results and in most cases are highly unoptimized. None of the approaches discussed thus far uses a formal and explicit representation of either the user context or the source-code structure. Whether the representation mechanism used is a relational database (Strathcona) or traditional graph-based representation (PARSEWeb, XSnippet and Prospector), we hypothesize that encoding the representation using ontology formalisms greatly improves the search and retrieval of relevant code snippets. Ontologies can naturally combine knowledge from multiple sources

(contexts) and then allow for computing entailments from this combined knowledge.

To explore this hypothesis, we use a purely semantic-based approach to source-code recommendation. This approach uses ontologies to represent the complex dependencies of code elements and allows for modeling a large class of inter-relationships and dependencies between various program elements.

3 Ontology representation of source code

An ontology is an explicit data model for a particular domain. It consists of machine-interpretable definitions of classes that formally describe domain-specific concepts, relationships between classes, structural properties of classes, and constraints, expressed as axioms. Due to their solid formal and reasoning foundation, ontologies have been successfully used by various research body to improve software engineering processes [15].

To explicitly represent the conceptual source-code knowledge of the user context as well as software libraries used in the user's project, we have created a Source-Code-Representation Ontology (referred to afterwards as SCRO). SCRO provides a base model for understanding the relationships and dependencies among source-code artifacts. It captures major concepts and features of object-oriented programs, including encapsulation, class and interface inheritance, method overloading, method overriding, and method signature information. SCRO's knowledge is represented using OWL-DL² ontology language. OWL is a web-based language used for capturing relationship semantics among domain concepts, OWL-DL is a subset of OWL based on Description Logic and has desirable computational properties for reasoning systems. OWL-DLs reasoning support allows for inferring additional knowledge and computing the classification hierarchy (subsumption reasoning).

SCRO defines various OWL classes and subclasses. These classes map directly to source-code elements and collectively represent the most important concepts found in object-oriented programs. Furthermore, we define various object properties, sub-properties, and ontological axioms within SCRO to represent the various relationships among ontological concepts. SCRO is precise, well documented, and designed with ontology reuse in mind. The availability of SCRO online [1], allows researchers to reuse or extend its representational components to support any semantic-based application that requires source code knowledge.

After having created the ontology structure, one next needs to populate the knowledge base with ontological instances (OWL individuals) that represent various concepts in the ontology. In the context of source-code recommendation, we need to populate SCRO with instances from various

sources. Remember that we are interested in tracing method call sequences to produce a code snippet for a particular object instantiation task. As the main source of information for tracing such sequences, we consider the framework(s) currently being reused. Furthermore, to rank the retrieved snippets, we also take into account the user context depicted by the current project under development.

To that end, we have built a knowledge generator for Java. The generated semantic instances are serialized using the RDF³ language. RDF is suitable for describing resources and provides a data model for representing machine-processable semantics of data. For each library or framework parsed, our knowledge generator builds an RDF ontology that conforms to SCROs descriptions of source code. This way, we provide a clean separation of the explicit OWL vocabulary with its associated schema definitions represented in SCRO from the metadata encoded in RDF. Furthermore, our knowledge generator is fully automatic and efficient, it took only 4.5 seconds to parse, process, and generate a knowledge base for the entire Jena framework. For an extended description of SCRO, the process of knowledge population, and samples of our knowledge extractor subsystem, we refer the reader to our ontologies website [1].

3.1 Recommendation scenario

Consider a second example of a developer trying to reuse the Jena framework. The developer wishes to programmatically construct a fragment of a RDF model based on a template in a given query. She would start with a `Query` object obtained from a `String` representation of the query and wishes to end up with a `Model` object representing the newly constructed RDF model. For a developer who is unfamiliar with this API, accomplishing this task may not be easy, in part because there are some intermediate steps and various static method invocations for instantiating the desired `Model` object; these steps are outlined in the following code snippet:

```
Query query =
    QueryFactory.create(queryString);
QueryExecution qe =
    QueryExecutionFactory.create(query);
Model m = qe.executeConstruct();
```

Approaches based on heuristic-based queries [5], first use a set of hard-coded heuristics to extract the structural description of the code being developed. These approaches then match this description against the structural descriptions found in a repository populated with sample client code. The result returned to the user is a set of methods retrieved from the repository after applying all the heuristics.

²<http://www.w3.org/TR/owl-guide/>

³<http://www.w3.org/TR/rdf-primer>

Ideally, the method embodies the code snippet shown above is included in this set. Notably, these kind of queries can be easily answered using our proposed approach: SCRO currently captures all structural descriptions of source code, ontologies can be populated through parsing, and heuristics can be applied to the generated knowledge using semantic rules. Ontologies, in this case, provide expandable solutions for extracting new knowledge when it becomes available, as well as flexible means for encoding heuristics and ranking the results.

However, we instead answer an object-instantiation query by constructing a code snippet. This snippet then shows how to get a handle to the `Model` object in the code snippet above. In our approach, we rely on an RDF-graph representation of the application framework being reused. RDF is a flexible and extendable data-representation and graph-based modeling language. Unlike traditional graph-representation mechanisms, RDF graphs provide a precise description of resources and are capable of encoding metadata in its nodes (represented by OWL classes) and edges (represented by OWL object properties). Thus, RDF graphs provide flexible means for semantic reasoning and deductive querying. Figure 1 shows a partial RDF-graph representation of the code snippet we discussed above. As described in Section 3, we can obtain this representation by parsing the Jena application framework. Nodes in the graph show instances of OWL classes that are specified in the knowledge base, edges, however, are SCROs object properties. For example, `hasOutputType` is a functional property that represents the method’s return type. `hasInputType` represents the type of a method’s formal parameter, the inverse of this property is `isInputTypeOf`. Inverse properties are included so we can traverse the RDF graph in both directions.

Given this directed RDF graph representation, we construct a code snippet through a guided brute-force graph-traversal search starting at the node that represents the source type (`Query`), enumerating all possible path candidates to the given destination type (`Model`). However, the cost of traversing and querying a large RDF model is susceptible to increase as the number of OWL individuals in the knowledge base increases. Combined with inference provided by the reasoner, this can be an expensive combination. To avoid reasoning overhead, we use the reasoner to obtain the inference closure of the original model. We then save the result which includes the computed entailments into a plain ontology model, with no inference engine attached to it. Therefore, the new model will be used for further processing. This way, we maintained the benefits of inference but avoided the added costs implied by using the reasoner. Furthermore, since we are only interested in object-instantiation queries of the form (*Source object* \Rightarrow *Destination object*), not every path in the graph is of in-

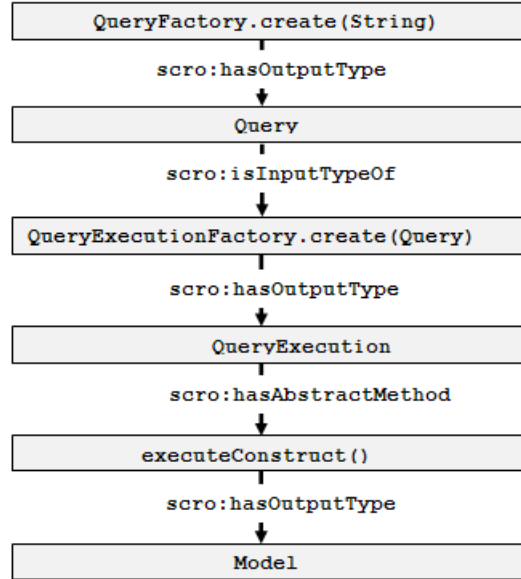


Figure 1. Answering: $Query \Rightarrow Model$

terest to us. Therefore, we further restrict the obtained plain model to only those RDF nodes and edges that can ever be used in a path that represents a code snippet from the source object to the corresponding destination object. The obtained model will be used as the basis for graph traversal, querying, and hence snippet construction. We call this model the *Snippet Model*.

4 Handling Downcasts

The scenario discussed in Section 3.1 only shows features that one can obtain from parsing a given framework’s API signatures. In particular, some complex application frameworks do not usually provide type-specific methods that may free the developer from using narrowing reference conversions (downcasts) to complete a particular programming task. At its current state, the *Snippet Model* and its associated RDF graph has no support for such features. We thus introduce the `hasActualOutputType` OWL property from SCRO. This property represents the actual runtime type of a given method. Therefore, we enrich the RDF graph with an edge labeled with `hasActualOutputType` that leaves a given method and enters every possible runtime type of this method. When we construct a code snippet, `hasActualOutputType` is treated as an expression casting the result of the method down to its actual return type.

In order to obtain the precise return type of API methods we rely on inter-procedural points-to and call graph analysis. Points-to analysis [3] is a static program analysis technique that analyzes a sample program in order to obtain

precise reference and call-target information. Among other usage areas, points-to analysis techniques have been used in solving complex software engineering problems. The essence of using these techniques in our approach is to analyze each API method in order to obtain a points-to set of possible runtime return types of the target method. Once obtained, this information is used to enrich the snippet model with paths, along the `hasActualOutputType` property, between a given method and its actual runtime type. This information will be used to determine the legality of a type cast when search is performed at a later stage. In order to perform this pre-processing step, we use the Soot framework [14], a popular program analysis and optimization framework for Java programs. Soot is capable of performing inter-procedural data-flow analysis on a whole-program points-to graph mode. However, this precise analysis requires creating program entry points. Although we are currently working on methods that make points-to analysis viable on programs that have no distinguished entry point, this is still ongoing work. Therefore, we currently write appropriate entry points manually, by providing a main class with a main method that exercises the API in question.

Suggesting a code snippet that contains a downcast is not the norm. In fact, the need to use a downcast reveals the hidden complexity of the underlying framework. Consider a programmer coding for the Eclipse API who wishes to obtain a handle of `JavaInspectExpression` representing the currently selected expression in the debugger represented by an object of type `IDebugView`. The sample code below details a sample solution:

```
IDebugView debugger = ...
Viewer viewer = debugger.getViewer();
IStructuredSelection sel =
    (IStructuredSelection)
        viewer.getSelection();
JavaInspectExpression expr =
    (JavaInspectExpression)
        sel.getFirstElement();
```

This example is in fact taken from the Prospector tool as a motivational example for using a corpus of client code to synthesize code snippets with downcasts. In such rare situations where the RDF graph will not find a solution path, we rely again on points-to-analysis to construct such snippet. Figure 2(a) shows part of the RDF graph where the dead end is reached. `ISelection` has in fact a single method of no relevance, `isEmpty()`⁴ and there is no way for the traversal algorithm to proceed further. We thus infer the runtime return type of method `getSelection()` by using flow-insensitive points-to analysis on the partially

⁴Methods that return a value of primitive type are excluded from the snippet graph

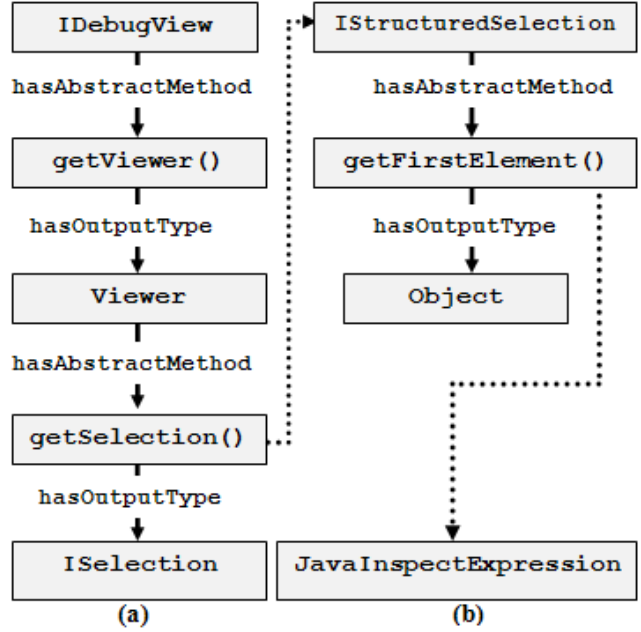


Figure 2. `IDebugView` \Rightarrow `JavaInspectExpression`

constructed snippet. Upon completion, this analysis concludes that `IStructuredSelection` is the actual return type of `getSelection()` at that particular call site (dashed lines represent `hasActualOutputType`). The same process is repeated for `getFirstElement()` as shown in Fig. 2(b). We thus avoided using a static corpus of client code that may in fact have no sample for answering this query.

5 Snippet ranking and selection

Since every obtained code snippet represent a path in the graph, it is notable that answering a given query may result in a large number of paths, each of which is a solution candidate. However, not all solutions of the same degree of relevancy to the user. It is also notable that no recommendation mechanism can rely on the type system to identify best candidate solutions, therefore, heuristics can be used to rank the results based on relevancy to the task at hand. In our approach, we use the path size heuristic as well as user context heuristics.

Shortest-Path-First (SPF) is simple yet proven effective heuristic. It assigns top rank to the shortest path in the graph. This heuristic is a variation of the code length heuristic proposed by Prospector[7] and used by others. However, in our case, a path size represents the number of RDF statements that are necessary to compose the snippet. Simply, a RDF statement is an assertion that a relation holds between two resources. For instance, the statement (`executeConstruct()` `hasOutputType Model`) asserts that the relation `ha-`

`sOutputType` holds between `executeConstruct()` and `Model`. Therefore, the size of the path in Fig. 1 is five.

While SPF clusters and ranks paths based on their size, context-based heuristic assigns higher ranks to paths that better fit within the current user context. We thus analyze the code that is currently being developed by the user, then we create a context profile that include all visible types that are either declared by the user or inherited in the user’s context. We further analyze each retrieved path in terms of the new types that this path will introduce into the current context. For example, a particular method invocation may have an argument that requires instantiating and thus introducing a new type into the current context. Naturally, code snippets that introduce more types should be assigned a lower rank value. However, if a newly introduced type is found in the context profile, it will not count against the enclosing path. This is entirely based on a simple scoring procedure that accounts for the number of newly introduced types and their visibility in the context profile. These heuristics are simple, easy to implement, and helped improve our results.

6 Implementation and evaluation

In order to investigate and evaluate the value of ontologies and points-to analysis in source code recommendation, we have implemented RECOs, a prototype object instantiation and recommendation system. RECOs is currently combined with a tool we have developed for detecting design pattern instances in object-oriented frameworks [2]. This combination is meant to promote multiple levels of software understanding and knowledge reuse. It is evident, based on empirical studies [6], that code examples are undoubtedly necessary for understanding framework usage, however, examples alone may not be enough to achieve higher potential of systematic software reuse. To be truly effective, developers need to learn the design knowledge implemented in these reusable frameworks. Our current implementation provides both advantages in one tool.

In order to use RECOs, one needs to provide the location of the framework’s binaries. The knowledge extractor subsystem automatically parses the jar files and generates a RDF ontology representing the structural description of the framework’s API. This ontology is classified by the reasoner to generate semantic entailments and ensure proper conformance to SCRO’s vocabulary and constraints. Once classified, the system automatically generates the snippet model that is subsequently enriched with ontological instances obtained via points-to analysis as described in Sections 3.1 and 4, respectively. The final ontology serves as the basis for answering object instantiation queries for that particular framework. In fact, this snippet ontology need not be changed until a new framework release is available. The recommender subsystem accepts a user query (using a

very simple input form for entering the source and destination objects), performs graph traversal, selects and ranks appropriate paths, and generates a custom code snippet for each path. This subsystem is independent from the knowledge generator subsystem. A user need only to configure it with the location of the snippet ontology, the directory of the code currently being edited, and the number of code snippets returned when a query is executed (default is 15).

In order to assess the benefits acquired by our approach, we conducted multiple experiments and case studies. The fundamental guiding hypotheses we test in these experiments are:

H1 *Ontology-based representation of source code knowledge improves search precision, and provides better recommendation of code snippets.*

H2 *Interprocedural points-to analysis techniques relieve a recommendation system from relying on a repository populated with sample client code.*

H3 *Contextual information provides better ranking and filtering of recommended items.*

6.1 Case study: framework usage

This experiment is designed to evaluate RECOs accuracy for answering object instantiation specific tasks when reusing an application framework. Jena has been selected for testing RECOs since it has an informative mailing list and forum for Jena developers⁵. This allows us to collect actual coding problems posted by developers and answered by the Jena support team. Furthermore, our familiarity with Jena allows us to inform the study with other coding tasks that are designed to test different aspects of our approach. Therefore, we have selected thirteen tasks that vary in their complexity, ranging from a simple constructor or static method call to a complex sequence of expressions. Table 1 shows statistics about these tasks after being expressed as object instantiation queries.

For each task, an environment has been setup such that the desired code snippet is left incomplete. We then instructed RECOs to fill in the missing code. We remind the reader that a desired solution to a given query may not be completely ready for immediate insertion in user code. In some cases, the user still need to issue another query to instantiate one or more objects introduced by the solution (e.g., an argument in a method call or an intermediate object within the sequence). Although not complete, such solution is still desired as long as it contains a valid sequence that would ultimately complete the task in the given context. Consider task T3 for example, the retrieved code snippet requires an object of type `OntModel` to be present

⁵<http://tech.groups.yahoo.com/group/jena-dev/>

Table 1. RECOS framework usage results

Task	Source	Destination	Context	Rk. ¹
T1	Statement	RDFDataType	-	1
T2	RDFList	IntersectionClass	-	1
T3	Resource	ComplementClass	-	3
T4	<i>Null</i>	OntModel	String	1
T5	StmtIterator	RDFNode	Property	2
T6	Prologue	Query	-	1
T7	Property	Selector	RDFNode	1
T8	ResultBinding	RDFNode	-	1
T9	Statement	RDFList	Property, OntClass	2
T10	IDBConnection	ModelRDB	-	1
T11	String	Individual	OntModel, String	0
T12	Query	Model	String, OntModel	3
T13	Resource	OntModel	Model	3

¹ Rk: Rank of desired solution if found, 0 otherwise.

for the code to compile. RECOS generates an intermediate variable that suggests a need to instantiate this object if not already instantiated. Typically, these objects can be instantiated with a query specifying only the destination object as seen in task T4. RECOS accepts queries of the form *Null* \Rightarrow *Destination*, this is useful for creating objects using a simple constructor or static method call, or when the source object is completely unknown. Usually, generalized queries of this form produce many hits. In fact, since we only rely on API signatures, the number of hits is, in many cases, large. However, we believe that the precise ontology-based representation of API code combined with heuristics produces a good rank of the desired result. This shows a clear support for hypotheses H1 and H2. Hypothesis H2 is also supported in part by tasks T8-T10. In task T8 for example, an intermediate method that returns an object of type `Object` needed to be converted to `RDFNode`. Points-to analysis inferred that this cast is possible, thus, avoided long and undesired paths.

In some case, it was extremely difficult to infer the user’s intent. Consider task T11 for example, RECOS returns many hits that do not complete the task in question. This task in fact shed a light on one of the difficulties faced by recommendation systems in general. Dealing with the highly polymorphic `String` objects affects precision and requires more sophisticated approaches to be handled properly. In this example, it was not clear to RECOS whether the user wants to get an individual represented by a `String` URI from an ontology model or in fact she wants to create an individual representation from an ontology class. On the other hand, API methods that are heavily overloaded have their own affects on ranking results. In task T12, RECOS ranked the desired solution third and generates many hits. This snippet as described in Section 3.1 requires invoking a heavily overloaded method that creates

a `QueryExecution` to execute over the ontology model. However, in this case, context heuristics filtered out plenty of paths that would otherwise get a higher rank. Hypothesis H3 was also supported by T13. RECOS, in this case, filtered out many paths that would have introduced a new object (e.g., `OntModelSpec`) to the user context. With the exception of the `XSnippet` tool, we do not expect other tools to perform well when context is necessary. However, `XSnippet` is tied to querying specific frameworks and the context used depends on the structure of the examples found in the repository of sample code. This indicates that `XSnippet` has to traverse the graph representation of the repository every time the user issues a query.

6.2 Comparison with other approaches

It is not trivial to compare code search tools due to the obvious lack of standard benchmarks and tool availability issues⁶. We thus extend a case study proposed by Thummalapenta and Xie [13] and used to evaluate `PARSEWeb` against `Prospector` and `Strathcona` search tools. This case study is based on the `Logic` example project of the `Eclipse Graphical Editing Framework(GEF)`⁷. The authors proposed ten programming tasks, shown in Table 2, expressed as object instantiation queries. We parsed the needed jar files, generated the `GEF Snippet Model`, and instructed `RECOS` to find solutions for each task.

Table 2. GEF Logic tasks & results

Task	Source	Destination	Rank
T1	<code>IPageSite</code>	<code>IActionBars</code>	1
T2	<code>ActionRegistry</code>	<code>IAction</code>	1
T3	<code>ActionRegistry</code>	<code>ContextMenuProvider</code>	1
T4	<code>IPageSite</code>	<code>ISelectionProvider</code>	1
T5	<code>IPageSite</code>	<code>IToolBarManager</code>	1
T6	<code>String</code>	<code>ImageDescriptor</code>	1
T7	<code>Composite</code>	<code>Control</code>	10
T8	<code>Composite</code>	<code>Canvas</code>	7
T9	<code>GraphicalViewerThumbnail</code>	<code>Scrollable</code>	0
T10	<code>GraphicalViewer</code>	<code>IFigure</code>	0

As observed in Table 2, `RECOS` found solutions for all tasks except T9 and T10. T9 was in fact unfeasible since we could not verify the existence of the `Source` object in the library. T10 was not answered by any of the tools. In fact, `PARSEWeb` was unable to find a solution for T3, `Prospector` and `Strathcona` could not answer T6-T8. Consider T8 for example, the shortest solution would be to invoke an existing `Canvas` constructor that accepts

⁶Closely related tools, `XSnippet` and `Prospector` are currently unavailable. `PARSEWeb` is not currently communicating with `Google Code Search(GCS)` since `GCS` has changed its interface

⁷<http://www.eclipse.org/gef/>

`Composite` in its parameter list. However, a desired solution based on the user’s context is, in fact, to instantiate `PageBook`, a sub-type of `Composite`, and pass that object to the constructor. Without our context heuristics, `RE-COS` would have ranked this answer further down in the list. `PARSEWeb` ranked this higher since it utilizes the usage frequency heuristic. `Prospector`, on the other hand, could not answer the query or perhaps the answer did not show up in the list⁸ because its ranking mechanism is based mostly on the length heuristic.

These results show a clear support for our three hypotheses. However, `PARSEWeb` outperforms `RECOs` and `Prospector`⁹ in the number of retrieved results. This is expected due to our reliance on API structures, but the effect of the final results was greatly reduced due to the nature of the semantic representation and organization of API knowledge. Furthermore, `PARSEWeb` ranked some of the desired solutions higher. `PARSEWeb` performs sequence post-processing and clustering that appears to improve the total number of retrieved results and plays a role in ranking similar sequences. However, `PARSEWeb` relies on incomplete code fragments obtained from `GCS` and has no access to API information. Therefore, it must use various heuristics for type resolution. These heuristics, however, may not work when the downloaded code contains a complex sequence of method calls that was not used in initialization expressions. Achieving perfection in code search is near impossible, however, we believe that `RECOs` internal mechanisms proved effective, and in the majority of cases, show a clear support for our three hypotheses.

7 Discussion and future work

Ontologies have been widely recognized as effective means for knowledge representation. On the other hand, points-to analysis techniques provide effective mechanisms for type inference. Both techniques have been independently used in solving different software engineering problems. In this paper, we proposed an approach that combines the strength of both techniques to improve search for relevant source-code snippets. We have also developed `RECOs`, a code search tool for object instantiation specific queries. `RECOs` is currently not tied to a particular IDE. However, we are currently integrating `RECOs` into `Eclipse` as part of a comprehensive tool for program understanding and knowledge reuse. In addition to snippet recommendation and design recovery, this tool will be used to recommend reusable components (e.g., finding source and destination objects). We are also investigating the application of

semantic annotations and domain ontologies to improving search precision and ranking.

References

- [1] A. Alnusair and T. Zhao. Source Code Ontology (SCRO) and examples of automatic ontology population. <http://www.cs.uwm.edu/~alnusair/ontologies>.
- [2] A. Alnusair and T. Zhao. Towards a model-driven approach for reverse engineering design patterns. In *2nd International Workshop on Transforming and Weaving Ontologies in Model Driven Engineering (TWO-MDE’09)*, 2009.
- [3] M. Emami, G. Rakesh, and L. Hendren. Context-sensitive interprocedural points-to analysis in the presence of function pointers. In *ACM Conference on Programming Language Design and Implementation (PLDI)*, pages 242–256, 1994.
- [4] T. R. Gruber. A translation approach to portable ontology specification. *Knowledge Acquisition*, 5(2):192–220, 1993.
- [5] R. Holmes and G. C. Murphy. Using structural context to recommend source code examples. In *International Conference on Software Engineering (ICSE)*, pages 117–125, 2005.
- [6] D. Hou. Investigating the effects of framework design knowledge in example-based framework learning. In *IEEE International Conference on Software Maintenance*, pages 37–46, 2008.
- [7] D. Mandelin, L. Xu, L. Bodik, and D. Kimelman. Jungloid mining: helping to navigate the API jungle. In *Programming language design and implementation*, pages 48–61, 2005.
- [8] A. Michail. Data mining library reuse patterns using generalized association rules. In *International Conference on Software Engineering (ICSE)*, pages 167–176, 2000.
- [9] D. Poshyvanyk, A. Marcus, and Y. Dong. JIRiSS—an eclipse plug-in for source code exploration. In *IEEE Conference on Program Comprehension*, pages 252–255, 2006.
- [10] M. P. Robillard, R. J. Walker, and T. Zimmermann. Recommendation systems for software engineering. *IEEE Software*, 2010.
- [11] R. Sindhgatta. Using an information retrieval system to retrieve source code samples. In *International Conference on Software Engineering*, pages 905–908, 2006.
- [12] N. Tansalarak and K. Claypool. XSnippet: mining for sample code. In *Conference on Object-Oriented Programming, Systems, Languages, and Applications*, pages 413–430, 2006.
- [13] S. Thummalapenta and T. Xie. `PARSEWeb`: a programmer assistant for reusing open source code on the web. In *IEEE/ACM International Conference on Automated Software Engineering*, pages 204–213, 2007.
- [14] R. Vallée-Rai, P. Co, E. Gagnon, L. Hendren, P. Lam, and V. Sundaresan. Soot - a Java bytecode optimization framework. In *Conference of the Centre for Advanced Studies on Collaborative Research*, pages 242–256, 1999.
- [15] Y. Zhao, J. Dong, and T. Peng. Ontology classification for semantic web based software engineering. *IEEE Transactions on Services Computing*, 2(4):303–317, 2009.
- [16] H. Zhong, T. Xie, L. Zhang, J. Pei, and H. Mei. MAPO: Mining and recommending api usage patterns. In *European Conference on Object Oriented Programming (ECOOP’09)*, pages 318–343, 2009.

⁸Prospector was configured to show only the first 12 results

⁹Strathcona did not perform well in this experiment since it does not have a clear support for object instantiation queries