

Class Schema Evolution for Persistent Object-Oriented Software: Model, Empirical Study, and Automated Support

Marco Piccioni, *Member, IEEE*, Manuel Oriol, and Bertrand Meyer

Abstract—With the wide support for object serialization in object-oriented programming languages, persistent objects have become commonplace and most large object-oriented software systems rely on extensive amounts of persistent data. Such systems also evolve over time. Retrieving previously persisted objects from classes whose schema has changed is, however, difficult, and may lead to invalidating the consistency of the application. The ESCHER framework addresses these issues through an IDE-integrated approach that handles class schema evolution by managing versions of the code and generating transformation functions automatically. The infrastructure also enforces class invariants to prevent the introduction of potentially corrupt objects. This paper describes a model for class attribute changes, a measure for class evolution robustness, four empirical studies, and the design and implementation of the ESCHER system.

Index Terms—Versioning, persistence, serialization, object-oriented class schema evolution, IDE integration



1 INTRODUCTION

BETWEEN the time an object-oriented program writes objects to persistent storage and when another execution retrieves them, the program may have evolved; in particular, classes describing these objects may have changed. The need to fit old objects into new classes is the problem of schema evolution for object-oriented software. While many solutions have been suggested and are currently used, none has been to our knowledge generally accepted.

The techniques used to cope in practice rely heavily on manual effort by the developers, who must examine previous class versions and provide conversion code. This approach is not only tedious, but a threat to software reliability as it usually relies on tolerant retrieval algorithms that make questionable decisions about the key issue: how to avoid accepting semantically inconsistent objects into the retrieving system.

The purpose of the present work is to lay a solid foundation for a general, stable solution to the problem of object-oriented class schema evolution. The contribution includes four major components: empirical studies, a formal model, resulting both from the empirical studies and from object-oriented theory, a measure of persistence evolution robustness, and a schema evolution tool implementing the model.

Section 2 presents a motivating example and the corresponding ESCHER tool support. Section 3 describes previous

approaches to class schema evolution. Section 4 defines a formal model representing atomic attribute changes through a set of transformation functions, and suggests a definition of a measure for evolution robustness. The empirical studies in Section 5 are essential to understanding what kind of schema evolution actually happens in practice; without empirical evidence, proposed solutions to the schema evolution problem are bound to be off the mark. The first study analyzes the evolution of serializable classes in five versions of the `java.util` package. It highlights the acuteness of the problem by showing that 17.7 percent of the changes affect the persistence of the serializable classes, and justifies the set of atomic attribute changes considered in the model. The measure for evolution robustness is also calculated to provide further evidence. The second study covers, for balance, the history of five versions of the Eiffel base library classes. The findings confirm those of the first study and lead to a refined set of atomic attribute changes. Two more studies analyzing persistent classes in two well-known software projects, the Apache Tomcat web server and the EiffelStudio IDE (development environment), lead to similar conclusions. The implementation (Section 6) is the ESCHER tool, seamlessly integrated into EiffelStudio and supporting developers through the entire schema evolution process, from class release time to retrieval time. It includes version handling, code template generation for transformation functions, and a retrieval algorithm preventing the acceptance of inconsistent objects into the system. Section 7 evaluates ESCHER by analyzing which code changes ESCHER actually detects in the four projects mentioned above. It then analyzes the advantages and limitations of the contribution, drawing conclusions and suggesting future work. This paper extends previous work of Piccioni et al. [1] by adding three new studies (on the EiffelBase library, on the Java Apache Tomcat repository, and on the EiffelStudio IDE repository), by updating the model, by defining and applying a measure

- M. Piccioni and B. Meyer are with the Swiss Federal Institute of Technology (ETH), Clausiusstrasse 59, CH-8092 Zürich, Switzerland. E-mail: {marco.piccioni, Bertrand.Meyer}@inf.ethz.ch.
- M. Oriol is with ABB Switzerland Ltd Corporate Research, Segelhofstrasse 1K, CH-5405 Baden-Daettwil, Switzerland. E-mail: manuel.oriol@ch.abb.com.

Manuscript received 24 Feb. 2011; revised 20 Nov. 2011; accepted 29 Nov. 2011; published online 28 Dec. 2011.

Recommended for acceptance by M. Robillard.

For information on obtaining reprints of this article, please send e-mail to: tse@computer.org, and reference IEEECS Log Number TSE-2011-02-0056. Digital Object Identifier no. 10.1109/TSE.2011.123.

```

class
  BANK_ACCOUNT
create make

feature -- Initialization
  make
  -- Create a bank account with an initial deposit.
  do
    tot_deposits := 1
  end
feature -- Status report
  balance: INTEGER
  -- Account balance.
  do
    Result := tot_deposits - tot_withdrawals
  end
  info: STRING
  -- Some numeric information.
feature -- Basic operations
  deposit (sum: INTEGER)
  -- Add 'sum' to account.
  require
    sum_non_negative: sum > 0
  do
    tot_deposits := tot_deposits + sum
  ensure
    balance_correct : balance = old balance + sum
  end

  withdraw (sum: INTEGER)
  -- Retrieve 'sum' from account.
  require
    sum_non_negative: sum > 0
    has_sufficient_funds : sum < balance
  do
    tot_withdrawals := tot_withdrawals + sum
  ensure
    balance_correct : balance = old balance - sum
  end
feature {NONE} -- Implementation
  tot_deposits : INTEGER
  -- Total amount deposited.
  tot_withdrawals : INTEGER
  -- Total amount withdrawn.
invariant
  valid_account : tot_deposits > tot_withdrawals
end

```

Fig. 1. *BANK_ACCOUNT*, version 1.

for evolution robustness to all the classes in the `java.util` package, and by refining the tool implementation.

2 A SCHEMA EVOLUTION EXAMPLE AND ITS HANDLING IN ESCHER

This section presents an example scenario and shows how the ESCHER tool copes with it.

2.1 A Case of Object-Oriented Schema Evolution

Assume there is a software system storing objects of type *BANK_ACCOUNT* (see Fig. 1). The balance is computed on-demand from the total amount of deposits and withdrawals, and there is an explicit class invariant capturing the bank account intended semantics. This invariant is checked after the invocation of the constructor *make* of a bank account object, and also before and after the invocation of any other routine in the class. Version 1 then evolves into version 2 (Fig. 2), in which the query *balance* becomes an attribute updated every time a deposit or a withdraw takes place, and the attribute *info* becomes an integer. While in languages with a C-derived syntax like

```

class
  BANK_ACCOUNT
create make

feature -- Initialization
  make
  -- Create a bank account with an initial deposit.
  do
    balance := 1
  end
feature -- Status report
  balance: INTEGER
  -- Account balance.
  info: INTEGER
  -- Some numeric information.
feature -- Basic operations
  deposit (sum: INTEGER)
  -- Add 'sum' to account.
  require
    sum_non_negative: sum > 0
  do
    balance := balance + sum
  ensure
    balance_correct : balance = old balance + sum
  end

  withdraw (sum: INTEGER)
  -- Retrieve 'sum' from account.
  require
    sum_non_negative: sum > 0
    has_sufficient_funds : sum < balance
  do
    balance := balance - sum
  ensure
    balance_correct : balance = old balance - sum
  end
invariant
  valid_account : balance > 0
end

```

Fig. 2. *BANK_ACCOUNT*, version 2.

Java we would need to modify client code to accommodate the fact that a method becomes an attribute, here clients are not affected because in Eiffel both attributes and argument-less functions can be accessed in the same way from outside the class, providing uniform access to the class itself. Note also how the invariant is now expressed in terms of the new attribute *balance*. What happens when trying to retrieve an object stored with version 1 of class *BANK_ACCOUNT* into an object of version 2 depends on the retrieval system. Silently accepting objects having the attribute *balance* initialized to the default is in this case wrong. To understand why, imagine developers forget to code an appropriate transformation function assigning the difference between *tot_deposits* and *tot_withdrawals* stored in version 1 to the newly created *balance* attribute. Then *balance* will be zero for all retrieved objects, which is hardly ideal. If we consider the Eiffel version 2 of class *BANK_ACCOUNT* instead, at deserialization time there will be an invariant failure triggering an exception because the value of *balance* is not positive as the class invariant prescribes.

2.2 The ESCHER Solution

ESCHER aims at supporting developers in evolving object-oriented applications. It integrates seamlessly into an IDE, keeping track of all class versions and system releases. When a coherent set of classes is released, developers are supported in writing transformation functions to migrate from one version to another. Depending on the specific

evolutionary changes automatically detected, ESCHER support can vary from a function template containing some standard initializations to a full-blown function body generation, where no additional developer intervention is needed. Focusing on the previous example, we now show a worst-case and a best-case scenario with respect to code generation. For the *balance* attribute, the tool recognizes that two attributes have disappeared and a new one has been added, but obviously needs the developer's help to figure out how to initialize it. On the bright side, a default value initialization will raise an exception because it violates the version 2 class invariant, so there is no risk of accepting inconsistent objects into the system. For the *info* attribute, ESCHER recognizes the type change and provides a full code generation of the transformation function body. Finally, ESCHER includes a deserialization algorithm that is aware of the presence (or absence) of transformation functions between any two versions of each class, triggering deserialization failures if the required transformation function is not present, or if the deserialized object does not satisfy the new version class invariant.

3 PREVIOUS APPROACHES

The issues arising from schema evolution are widely acknowledged. They affect object-oriented databases, relational databases, and all programming languages providing a serialization facility.

3.1 The Two-Version Approach

Given a class whose objects we need to serialize, the most widespread approach takes into account the "current" version and the "stored" one. To convert objects between the two, developers typically implement a transformation function that should be able to handle all possible stored versions of the data.

In Java, a class can enable future serialization of its instances by implementing the `Serializable` interface or its descendant `Externalizable` [2], [3], [4]. Developers can provide custom deserialization methods to help establishing the class invariant.

The .NET framework, in a way similar to Java and starting from version 2.0, provides `Version Tolerant Serialization` (VTS) [5]. In spite of the name, VTS does not implement true version support. The serialized object and all the referenced objects in the object graph have to be tagged with the `Serializable` attribute. By using "attributes" instead of interfaces the mechanism is decoupled from the class hierarchy. As in Java, object constructors are not invoked during the deserialization process, so implicit class invariant violations may happen if developers do not provide custom deserialization methods.

Eiffel's serialization mechanism also offers a solution in which all conflicts are resolved in one class. Custom deserialization behavior is available by inheriting from a `MISMATCH_CORRECTOR` class and redefining a callback to re-establish the class invariant. As a difference with the solutions above, it is easier to detect class invariant violations because of the native support for Design by Contract [6].

In our view, the serialization mechanisms of O-O languages are converging toward a full-fledged solution

such as an Object Oriented Database Management System (OODBMS). To assess such a solution, we examine the db4o object oriented database [7].¹ Db4o stores objects as they are, without polluting them with persistence code. With respect to schema evolution handling, if developers need a custom behavior to establish the retrieving object class invariant, they have two possibilities: choose to use reflectively invoked methods in the object class or register listeners to specific `ObjectContainer` events outside the object class. As seen for Java and .NET, there is not explicit support for class invariant violation.

The Orion object-oriented database [8] introduces a model recognizing the importance of invariants to validate conversions. Compared to ESCHER, the Orion model is more complex because it deals with a broader set of code changes, mostly involving single and multiple inheritance, while ESCHER considers the object's flattened form. Though simpler, the ESCHER model is in one case more general because it allows switching to any type when detecting an attribute type change, while Orion only allows switching to an ancestor type. ESCHER also offers limited support for "attribute renamed," while Orion does not take any action in this case. Finally, in Orion invariants are mostly related to keeping the system consistent with respect to inheritance, while ESCHER uses class invariants to validate the conversions semantically.

3.2 The Multiversion Approach

By storing all versions of a class, the multiversion approach makes it easier to provide more specific schema evolution support. ENCORE [9] is an early example enforcing a serialized form limited to the attributes in the class interface, and proposing an ad hoc constraint language for interproperty constraints.

A general framework similar to ESCHER and applied to an OODB was proposed by Lautemann [10]. It required a schema designer role conceptually distinct from the application developer role, and envisioned a completely transparent schema evolution that we believe cannot be achieved.

A more recent example of a multiversion OODB is Versant Fastobjects.² It consists of two main storage areas: a dictionary, which carries the class definitions for objects in the database, and the actual database, containing the serialized objects forms. The dictionary is necessary to understand the structure of the objects stored in the database. The delicate issue consists of modifying the dictionary after having stored some objects. The dictionary keeps versions for every class, so when a class schema modification is detected, old objects are "transparently," lazily, and possibly inconsistently converted.

To convert a stored object of a certain version into an object of the current version, CLOSQL uses update or backdate routines [11]. An unfortunate requirement is that a database administrator is needed every time a class is created to specify which update or backdate routines have to be executed. As in our approach, the update and backdate routines for a certain version are kept all together, in this case in an "update method." CLOSQL supports

1. <http://developer.db4o.com/Home.aspx> (last visited: 20/11/11).

2. www.versant.com (last visited: 20/11/11).

“linear versioning”: A new version can only be generated from the latest version. In contrast to CLOSQL, we handle transformation functions also from older versions with respect to the last one.

GemStone is a computationally complete Smalltalk extension providing data definition, data manipulation, and query facilities to persist objects across executions [12]. Each class in GemStone has an associated class history containing all the previous versions. Developers decide when to define a new version and when to override ad hoc migration methods in the destination class.

The Oracle relational database, starting from version 11g release 2, has introduced “editions” to offer better support for online application upgrade [13]. The concept is more restrictive with respect to what the authors call releases because it is limited to “editionable” object types, which do not include, for example, tables and java classes.

When a certain class evolves over time, it may be considered a different type and named differently, or it may be considered the same type, keeping the same name but providing other means of taking into account the different inner structure and semantics. While the second option is mainstream, the first one has been previously explored by Piccioni et al. [14]. Every new version of a class is given a new name including version information. Versions can therefore be considered full citizens of the type system. To cope with this scenario, Eiffel provides `converters` [15]. Converters are a language mechanism intended to cover those cases in which inheritance is not appropriate, for example, to convert string implementations across different systems. A downside of this approach is the lack of scalability. In fact, every type should provide converters for each previous version, therefore clogging the corresponding class code with a potentially high number of transformation functions.

Another example of making different, parallel versions explicit can be found in `UpgradeJ` [16], though focused neither on object persistence nor on enforcing semantical consistency.

The automated detection of code changes has been extensively explored in the literature, both with respect to existing libraries and software configuration systems, for example by Dig et al. [17]. One interesting approach proposes that schema evolution can be solved by devising bidirectional transformations [18], [19]. Though acknowledging the advantages of generating the inverse transformations (mostly) for free from a bidirectional domain specific language, we believe that using it would be an excessive burden on the developer. We prefer to keep one programming language and provide transformation functions with explicit contracts to improve their reliability.

3.3 Other Approaches

Orthogonal Persistence Java (OPJ) is an extension of the Java Language Specification adding “orthogonal persistence” capabilities to the Java platform [20], [21]. The well-known prototype implementation for OPJ, PJama [22], [23], [24] is an extension of the Java Virtual machine together with a persistent store in which the state of an executing application is kept. The system state is check-pointed atomically and periodically to be able to recover from

exceptions and crashes. PJama provides an approach to schema evolution that involves persisting both objects and classes. To perform conversions between classes, PJama developers use a small API and a standalone, command-line utility. In case the changes are validated, objects are typically converted eagerly.

Several authors suggest assisting metamodel evolution by stepwise adaptation in a similar fashion with respect to what is presented in this paper, using a transformational approach and a classification of modifications [25], [26], [27], [28]. Analyzing the AST in metamodels and applying similarity metrics to detect changes has been done by Falleri et al. [29]. All the cited works about metamodels are of interest because class schemas are a particular kind of metamodels.

The automatic generation of transformation functions through type transformers in C programs has also been described by Neamtiu et al. [30]. It focuses on updating structs whose layout might evolve. However, it is not per se linked to object-orientation and it does not benefit from having a model or an integration into an IDE.

The formal object model proposed by the ODMG standard [31] encompasses more definitions than ESCHER’s, but it is missing some modifications on attributes that we have seen being relevant in practice, like renaming and changing type [32].

Schema evolution issues are also relevant for relational databases [33]. The inference rules we derive in Section 4 are related to those implemented in many schema matching tools available [34]. In particular, `TransScm` and `Tess` [35], [36] view a schema as a set of types and provide two-version transformations following certain rules. Human intervention is required to provide new rules or to select match candidates. Apart from the different transformation functions semantics and implementation, the most important differences are that ESCHER keeps track of all versions of a class schema and limits the conversions to those for which our data surveys show potential application.

The PRISM workbench [37], [38] presents some analogies with ESCHER, namely, the use of a GUI-based tool to facilitate schema evolution and a certain degree of automation in the process. The main difference we noticed is that in PRISM the database administrator has to manually input the Schema Modification Operators (SMOs), while in our case they are detected automatically while the system keeps track of all versions. An additional difference is that the SMOs, though more numerous than our atomic attribute changes, do not include “column type change,” an analog of our “attribute type changed.”

4 THE ESCHER MODEL OF OBJECT-ORIENTED SOFTWARE EVOLUTION

This section presents a model of updates that particularly emphasizes the generation of transformation functions—contrary to previous models [39], [40].

The upper part of Fig. 3 presents a simplified syntax of class definitions in Eiffel programs. We omit the declaration both of routines and constraints on generic parameters as they are typically not included in the serialized form. We also do not explicitly consider inheritance, both single and

| | |
|---|---|
| $C, D \in Names$ $N, M \in GNames$ $name \in names$ $att ::= name : type$ $\quad name : N$ $type ::= C$ $\quad type[type]$ $\quad \mathbf{attached} C$ $\quad \mathbf{detachable} C$ $class ::= \mathbf{Class} C \mathbf{feature} att_1, \dots, att_n \mathbf{end}$ $\quad \mathbf{Class} C[N_1, \dots, N_i] \mathbf{feature} att_1, \dots, att_n \mathbf{end}$ | <i>class names</i> <i>names for generic parameters</i> <i>names for attributes</i> <i>attribute definition</i> <i>attribute of generic type</i> <i>class type</i> <i>generic derivation</i> <i>void – safe type</i> <i>non void – safe type</i> <i>class</i> <i>generic class</i> |
|---|---|

$$att_i \in class_0; A_{noChange(att_i)}(class_0) = class_0$$

$$att \notin \{att_1, \dots, att_n\}; A_{new(att)}(\mathbf{Class} C \dots \mathbf{feature} att_1, \dots, att_n \mathbf{end}) =$$

$$\mathbf{Class} C \dots \mathbf{feature} att_1, \dots, att_n, att \mathbf{end}$$

$$(att_i = name : N \wedge att'_i = name' : N) \vee (att_i = name : type \wedge att'_i = name' : type)$$

$$att_i \in class_0; class_0 = \mathbf{Class} C \dots \mathbf{feature} att_0, \dots, att_n \mathbf{end}$$

$$class_1 = \mathbf{Class} C \dots \mathbf{feature} att_0, \dots, att_{i-1}, att'_i, att_{i+1}, \dots, att_n \mathbf{end}$$

$$A_{nameChange(name, name')}(class_0) = class_1$$

$$type \neq type'; att_i = name : type \wedge att'_i = name : type'$$

$$att_i \in class_0; class_0 = \mathbf{Class} C \dots \mathbf{feature} att_0, \dots, att_n \mathbf{end}$$

$$class_1 = \mathbf{Class} C \dots \mathbf{feature} att_0, \dots, att_{i-1}, att'_i, att_{i+1}, \dots, att_n \mathbf{end}$$

$$A_{typeChange(name:type, name:type')}(class_0) = class_1$$

$$att_i = name : type$$

$$att_i \in class_0; class_0 = \mathbf{Class} C \dots \mathbf{feature} att_0, \dots, att_n \mathbf{end}$$

$$class_1 = \mathbf{Class} C \dots \mathbf{feature} att_0, \dots, att_{i-1}, att_{i+1}, \dots, att_n \mathbf{end}$$

$$A_{removeAttribute(name)}(class_0) = class_1$$

$$att_i = name : C$$

$$att_i \in class_0; class_0 = \mathbf{Class} C \dots \mathbf{feature} att_0, \dots, att_n \mathbf{end}$$

$$class_1 = \mathbf{Class} C \dots \mathbf{feature} att_0, \dots, att_{i-1}, name : \mathbf{attached} C, att_{i+1}, \dots, att_n \mathbf{end}$$

$$A_{addAttach(att_i)}(class_0) = class_1$$

Fig. 3. Classes and SMOs.

multiple, because we have access to the flattened class schema. This is a valid assumption as serialized objects are de facto flattened as well.

The following sections define atomic attribute changes and class transformations, how to generate a type transformer from a class transformation, and present strategies to extract type transformers.

4.1 Classifying Software Evolution: SMOs

Adapting the definition from Curino et al. [37], [38] to an object-oriented context, we call Schema Modification Operator an incremental change to a class schema, along with the semantic changes possibly occurring to any instances conforming to that schema. More formally, an SMO is a function modifying at most one attribute

$$A : class \mapsto class.$$

We define six standard SMOs.

- Attribute not changed.
- Attribute added.
- Attribute renamed.
- Attribute type changed.
- Attribute removed.
- Attribute attached added.

The bottom part of Fig. 3 shows the definition of the six kinds of SMOs.

A class transformation T_{A_1, \dots, A_N} going from one version of a class to another can then be described by a list of SMOs A_1, \dots, A_N ($N \geq 1$):

$$T_{A_1, \dots, A_N} : class \mapsto class,$$

such that

$$T_{A_1, \dots, A_N}(class_0) = (A_n \circ \dots \circ A_1)(class_0).$$

Note that any modifications to the attributes of a class can be described by some class transformation—for example, by deleting all attributes whose names are not present in the class anymore and adding the new ones. Thus, class transformations are complete with respect to attribute modifications. This can be easily seen: Suppose we have two classes C1 and C2 with n attributes and m attributes, respectively, and assume the generic parameter list is preserved between C1 and C2. Therefore, there always exists a class transformation from C1 to C2 that first deletes all the n attributes from C1 and then adds all the m attributes to C2. While there is always a decomposition, using a straightforward algorithm might not produce the best results, so devising a heuristic is essential to using this model as the basis of an implementation.

| | | |
|----------------------|--|----------------------------------|
| $otname \in OTNames$ | | <i>object transformers names</i> |
| e | ::= $e.e \mid \mathbf{oldc} \mid name$ | <i>expressions</i> |
| <i>instruction</i> | ::= $\mathbf{Result.name}:=e \mid \mathbf{noop} \mid otname(e, e)$ | <i>instructions</i> |
| <i>instructions</i> | ::= $instruction \mid instruction; instructions$ | <i>list of instructions</i> |
| ot | ::= $otname(oldc : oldC) : C \text{ instructions}$ | <i>object transformer</i> |

| | |
|---|---|
| $\mathcal{G}[\cdot] : (class \mapsto class) \mapsto ot$ | |
| $\mathcal{G}[T_{A_1, \dots, A_n}] = otname(oldc : oldC) : C \mathcal{S}[T_{A_1, \dots, A_n}]$ | |
| $\mathcal{S}[\cdot] : (class \mapsto class) \mapsto instructions$ | |
| $\mathcal{S}[A_{noChange}(att_i)]$ | = $\mathbf{Result.name}_i := \mathbf{oldc.name}_i$ |
| $\mathcal{S}[A_{new}(att)]$ | = $\mathbf{Result.name} := \mathcal{I}[\cdot]$ |
| | <i>where $att = name : \dots$</i> |
| $\mathcal{S}[A_{nameChange}(name:type, name':type)]$ | = $\mathbf{Result.name}' := \mathbf{oldc.name}$ |
| $\mathcal{S}[A_{typeChange}(name:type, name:type')]$ | = $\mathbf{Result.name} := \mathbf{oldc.name}$ |
| | <i>if $type'$ is assignable to $type$</i> |
| | = $\mathbf{Result.name} := otname(\mathbf{oldc.name})$ |
| | <i>if there is an object transformer $otname$ transforming an object of type $type$ into an object of type $type'$</i> |
| | = $\mathcal{W}[name_0]$ <i>otherwise</i> |
| $\mathcal{S}[A_{removeAttribute}(name_0)]$ | = $\mathcal{W}[name_0]$ |
| $\mathcal{S}[A_{addAttach}(att_i)]$ | = $\mathbf{Result.name}_i := \mathbf{oldc.name}_i$ |
| $\mathcal{S}[T_{A_1, \dots, A_n}]$ | = $\mathcal{S}[A_1]; \mathcal{S}[T_{A_2, \dots, A_n}]$ |

Fig. 4. Object transformers.

4.2 Object Transformers

While SMOs allow an easy representation of the static transformations of a class, this is not sufficient to generate the object transformers that create an instance of the new class from a serialized instance of the old class. The main reason is that there is a need for an explicit default initialization of new attributes. To define such initializations, programmers actually need to input default initialization values. To represent that, we use the token $\mathcal{I}[\cdot]$ (for input), evaluating into the next value in a list of inputs. We also add an output generator $\mathcal{W}[\cdot]$ (for warning), warning programmers of the impossibility of generating a translation or a removal operation and evaluating into “noop.” The upper part of Fig. 4 presents a succinct syntax for object transformers. The generation of object transformers from a class transformation can be expressed as transformation functions, as shown in the bottom part of Fig. 4. The function $\mathcal{G}[\cdot]$ and generates the associated object transformer. Finally, the function $\mathcal{S}[\cdot]$ takes a class transformation and generates the actual object transformer code.

4.3 Heuristics for Class Schema Evolution

While we are currently considering the extraction of transformations based on actual changes made by programmers using the IDE, the current solution relies on statically comparing the Abstract Syntax Trees (ASTs) to detect SMOs. Static comparisons imply that in certain situations, like the “Attribute renamed” SMO below, the outcome of the tool can only suggest the possibility that the specific SMO has happened. The extraction relies on a set of heuristics successively applied.

1. An attribute that does not change name and declared type generates an “Attribute not changed” SMO. A possible change in the attribute semantics will be taken care of by the retrieving object class invariant.

2. An attribute that does not change name but changes type between two versions generates an “Attribute type changed” SMO. The framework is always able to detect this SMO and, depending on the type involved, to provide a complete code generation (as in the example in Section 2). A change in the attribute semantics is taken into account by the new type and by the retrieving object class invariant.
3. An attribute att removed from the old version is a candidate for generating an “Attribute renamed” SMO if in the new version there is a new attribute att' having the same type as att . By only comparing the two classes’ ASTs it is not possible to determine if a rename actually occurred because the new attribute could be semantically unrelated to the old one. In this case, we would have an “Attribute removed” followed by an “Attribute added.” The framework detects that this SMO might have occurred and generates a transformation function template warning about the possibility of a rename. In this case, the retrieving object class invariant is of the essence in establishing the correct object semantics. A better way to deal with the issue would be to check whether the variables are used by the same clients in the same context. As it requires a global analysis of the code, we considered it was too time consuming and could actually be captured in a better way by tracking user direct use of a renaming IDE functionality.
4. An attribute of the new version that does not appear in the old version generates an “Attribute added” SMO. The framework always detects this SMO, initializes the attribute to its default, and suggests checking the new class invariant with respect to the old one to avoid deserialization failures.
5. An attribute att removed from the old version generates an “Attribute removed” SMO if in the

new version there is no new attribute *att'* having the same type as *att*. The framework always detects this SMO. Once again, the class invariant of the retrieving object will take care of validating the new version semantics.

6. An attribute of the old version that has been augmented with the keyword *attached* in the new version generates an "Attribute attached added" SMO. The implementation of this SMO into the framework is ongoing work. The idea is to make sure the attribute is attached by explicitly creating an object in the transformation function body, attaching it to the attribute, and coding a postcondition stating that the attribute should be nonvoid. Here, both the transformation function postcondition and the retrieving object class invariant are crucial to validate possible semantics changes.

To detect the SMOs, we iterate through the new class attributes, search for a match with attributes in the old class, and create a corresponding heuristic. We then repeat the process starting from the old class in order to gather more information, for example to find all the attributes that were removed.

4.4 P-evolution-Robustness (PER): A Measure for Object-Oriented Software Evolution Robustness

It can be helpful to statically measure a software system with respect to its robustness to software evolution. Consider a class with m versions. Given that set of versions, there are $m(m-1)$ possible pairs of versions, each pair representing a transformation function from a version v_i to a version v_j ($i, j = 1, \dots, m, i \neq j$). We define a measure of class p-evolution-robustness as the ratio of the transitive closure of the binary relation "there is a transformation function between versions v_i and v_j " to the total number of transformation functions noted above. For instance, considering the two versions of class *BANK_ACCOUNT* introduced in Section 2, if someone defines only the forward transformation from version 1 to version 2, their PER would evaluate to 0.5: one defined transformation out of a possible two (one in each direction). As another example, consider having a class with five versions and four transformation functions defined, going from version v_i to version v_{i+1} ($i = 1, \dots, 4$). Because of the transitive closure, there will be six more transformation functions assumed because the existence of transformation functions between v_i and v_{i+1} and between v_{i+1} and v_{i+2} implies the one between v_i and v_{i+2} . Ten transformation functions are therefore assumed (out of 20), giving a PER of 0.5.

The PER measure can be easily restricted to a single class version v_i by computing the ratio of transformation functions defined between v_i and any other version to the total number of transformation functions that in this case is $2(m-1)$.

The PER measure can also be extended to a system release by computing the average of the PERs of the classes in the release.

In Section 5, we show how we compute PER for a software system that does not implement the multiversion approach, namely, the *java.util* package and its classes.

5 EMPIRICAL EVALUATION: THE ESCHER MODEL VERSUS THE PRACTICE OF OBJECT-ORIENTED SCHEMA EVOLUTION

To assess the relevance of our model, we analyze four sets of classes in two programming languages, Java and Eiffel. We focus on persistence-related code changes. For each language, we choose one set of classes from a widely used data structures and utility library, and one other from a well-known software project. Therefore, our investigation takes into account both libraries and applications. The tables shown capture the essential data. The complete set of data is available for download.³

5.1 Persistence-Related Code Changes in Java

Advani et al. [41] already evaluated code changes in 15 open source Java systems and showed that "rename field," "move field," "rename method," and "move method" account for approximately 66 percent of the total identified code changes. We can additionally observe that "rename field" and "move field" alone account for 32 percent of the total. While these results are calculated on all classes of the considered systems, there is no evidence that classes meant to be persistent would exhibit the same characteristics. It then seems reasonable to conjecture that class schema evolution in object-oriented applications is a relevant issue. In particular, the fact that class schema evolution involves attributes in one out of three cases suggests that providing support to handle the issue is a reasonable aim. To improve on the previous study, we focus on classes intended to produce objects that will be persisted. We therefore study a subset of likely to-be-persisted classes from the Java package *java.util* to check if they evolve in a similar manner to nonpersistent ones. Java code was a conscious choice because it is widely used and because classes that might have persistent instances are easily recognizable, as they are marked as *Serializable*.

5.1.1 Analyzing *java.util*

The package contains classes modeling collections, dates, currencies, and locales. We considered the 22 classes in the package directly implementing the *Serializable* interface and analyzed them manually across five versions of the language: 1.2.2, 1.3.1, 1.4.2, 5.0, and 6.0. We chose these versions because they can be considered significant steps in the evolution of the Java language, bearing important changes in both the APIs and their implementations.⁴ We also took into consideration 22 code changes types, seven of which are directly relevant to the serialization process and are shown in Table 1. "Attribute initialization value changed" is considered because it may influence the class invariant. "Attribute becoming a constant" is relevant because constants are not serialized, and therefore it will be interpreted as "Attribute removed." In a similar fashion, "Constant becoming an attribute" will be interpreted as "Attribute added." In line with what was discovered previously, Table 1 shows that the persistence-related code changes consist of 17.7 percent of the total (115 out of 649).

3. <http://tinyurl.com/ESCHER-data> (last visited: 20/11/11).

4. http://en.wikipedia.org/wiki/Java_version_history (last visited: 20/11/11).

TABLE 1
Changes Found across Five Versions of `java.util`

| | attribute added | attribute removed | attribute renamed | attribute type changed | attribute value changed | attribute becoming constant | constant becoming attribute | non persistence-related | all |
|-------------|-----------------|-------------------|-------------------|------------------------|-------------------------|-----------------------------|-----------------------------|-------------------------|-----|
| 1.2.2-1.3.1 | 4 | 0 | 0 | 0 | 1 | 0 | 0 | 11 | 16 |
| 1.3.1-1.4.2 | 10 | 13 | 1 | 1 | 0 | 0 | 1 | 99 | 125 |
| 1.4.2-5.0 | 29 | 6 | 0 | 13 | 2 | 4 | 0 | 270 | 324 |
| 5.0-6.0 | 14 | 9 | 4 | 1 | 0 | 2 | 0 | 154 | 184 |
| All | 57 | 28 | 5 | 15 | 3 | 6 | 1 | 534 | 649 |

Note also that “Attribute added” and “Attribute removed” together constitute 74 percent of all persistence-related code changes (85 out of 115).

Table 2 shows the distribution of code changes across different classes and across all five versions of the `java.util` package. The data suggest that persistence-related code changes are sufficiently widespread among classes.

This study shows that classes whose instances might be serialized change over time. Moreover, 17.7 percent of these changes directly impact the capability of classes to deserialize instances of their previous versions. This is aggravated by the fact that there may be a high number of stored objects that need to be updated.

In class `java.util.Calendar`, version 1.2, we found two version-related attributes, `currentSerialVersion` and `serialVersionOnStream`, in addition to the usual `serialVersionUID`. As explained in the code comments, these were used to solve the problem of keeping track of the different object versions, and had predefined values suggested. By using ESCHER and its integrated version support, creating these two

TABLE 3
P-Evolution-Robustness of `ArrayList`

| Version interval | Fw. tr. | Bw. tr. | Attribute changes |
|------------------|---------|---------|-------------------------|
| 1.2.2-1.3.1 | 1 | 1 | no changes |
| 1.2.2-1.4.2 | 1 | 0 | 1 added |
| 1.2.2-1.5.0 | 0 | 0 | 1 added, 1 type changed |
| 1.2.2-6.0 | 0 | 0 | 1 added, 2 type changed |
| 1.3.1-1.4.2 | 1 | 0 | 1 added |
| 1.3.1-1.5.0 | 0 | 0 | 1 added, 1 type changed |
| 1.3.1-6.0 | 0 | 0 | 1 added, 2 type changed |
| 1.4.2-5.0 | 0 | 0 | 1 type changed |
| 1.4.2-6.0 | 0 | 0 | 2 type changed |
| 5.0-6.0 | 0 | 0 | 1 type changed |
| | 3 | 1 | PER: 0.2 |

attributes would have proven unnecessary and the transformation function would have not been cluttering the class itself. In conclusion, while we cannot be sure that deserializing will be performed in a semantically inconsistent way, it is likely that it will create issues at some point.

5.1.2 The P-Evolution-Robustness of `java.util`

In spite of the fact that Java implements a two-version approach (as seen in Section 3), we are able to compute the PER for the `java.util` package. While the detailed data are available for download,⁵ in Table 3 we show how we performed the calculation for class `ArrayList`. Between versions 1.2.2 and 1.3.1 there are no changes, so both the forward and the backward transformations work. From version 1.2.2 to 1.4.2, there is an attribute added for which only the forward transformation works. The reverse transformation is not handled, and so deserialization of objects of version 1.4.2 from version 1.2.2 does not work because removing an attribute is an incompatible change according to the Java Serialization Specification [4]. In a similar way, and taking into account that an attribute type changed is also an incompatible change, we obtain for `ArrayList` a total of four existing transformations out of 20 possible across the five considered versions. The resulting PER value is therefore 0.2. In Table 4, we show the value of PER for every class in the `java.util`, and for the package itself.

The data suggest that the PER measure can be a simple way to highlight an often neglected aspect of software development, and can be easily used to assess the quality of a certain software repository and compare it to others.

5.1.3 Analyzing the Apache Tomcat Code Base

We choose to analyze Apache Tomcat as an example of a well-known and widely used application. A web server typically needs to store information about shared values that need to survive the stateless HTTP protocol interaction. The Tomcat open source project was initially released in 1999. It includes Catalina as a servlet container, Coyote as HTTP connector and Jasper as JSP engine. We analyze the evolution of persistent classes taking into account the earliest and the latest releases available in the open source repository at the time of writing (3.0 and 7.0, respectively).⁶ In Table 5 we show, for brevity, the classes in Tomcat’s code base in

TABLE 2
Changes Found in Five Versions of `java.util`, by Class

| Class | Code changes | Persistence-related |
|------------------------------|--------------|---------------------|
| <code>ArrayList</code> | 18 | 3 |
| <code>BitSet</code> | 42 | 8 |
| <code>Calendar</code> | 52 | 24 |
| <code>Currency</code> | 3 | 2 |
| <code>Date</code> | 22 | 9 |
| <code>EnumMap</code> | 1 | 0 |
| <code>EnumSet</code> | 1 | 0 |
| <code>EventObject</code> | 1 | 1 |
| <code>HashMap</code> | 101 | 11 |
| <code>HashSet</code> | 9 | 2 |
| <code>HashTable</code> | 41 | 5 |
| <code>IdentityHashMap</code> | 20 | 2 |
| <code>LinkedHashSet</code> | 5 | 1 |
| <code>LinkedList</code> | 50 | 1 |
| <code>Locale</code> | 34 | 16 |
| <code>PriorityQueue</code> | 17 | 0 |
| <code>Random</code> | 13 | 7 |
| <code>TimeZone</code> | 28 | 7 |
| <code>TreeMap</code> | 122 | 13 |
| <code>TreeSet</code> | 39 | 3 |
| <code>UUID</code> | 0 | 0 |
| <code>Vector</code> | 30 | 0 |
| All | 649 | 115 |

5. <http://tinyurl.com/ESCHER-data> (last visited: 20/11/11).

6. <http://svn.apache.org/viewvc/tomcat/> (last visited: 20/11/11).

TABLE 4
P-Evolution-Robustness of `java.util`

| Class | Class PER |
|-----------------|-----------|
| ArrayList | 0.20 |
| BitSet | 0.40 |
| Calendar | 0.25 |
| Currency | 0.70 |
| Date | 0.30 |
| EnumMap | 1.00 |
| EnumSet | 1.00 |
| EventObject | 0.70 |
| HashMap | 0.10 |
| HashSet | 0.30 |
| HashTable | 0.30 |
| IdentityHashMap | 0.40 |
| LinkedHashSet | 0.40 |
| LinkedList | 0.40 |
| Locale | 0.20 |
| PriorityQueue | 1.00 |
| Random | 0.15 |
| TimeZone | 0.30 |
| TreeMap | 0.20 |
| TreeSet | 0.35 |
| UUID | 1.00 |
| Vector | 1.00 |
| Package PER | 0.48 |

which we actually found at least one SMO. This happened in 12 out of the 27 classes analyzed. The data confirm what we found in the case of the `java.util` library classes. We did not notice significantly more changes as we might have expected from an application as opposed to a library project. This in our opinion might be because every software project values stability and backward compatibility and tries to limit changes to the bare minimum.

5.2 Persistence-Related Code Changes in Eiffel

5.2.1 Eiffel Libraries Code Changes

We analyze 22 Eiffel classes, most of them direct counterparts of the ones previously analyzed in Java. In the five cases in which it is not possible to find a direct counterpart, we choose classes whose instances are likely to be serialized, like strings and data structures. We also analyze five releases, but the

TABLE 5
SMOs Found in Tomcat's Repository

| | attribute added | attribute removed | attribute renamed | attribute type changed | all |
|-----------------------|-----------------|-------------------|-------------------|------------------------|-----|
| ValueReference | 0 | 2 | 0 | 0 | 2 |
| ApplicationParameter | 0 | 0 | 0 | 1 | 1 |
| FilterDef | 1 | 0 | 0 | 0 | 1 |
| FilterMap | 2 | 0 | 2 | 2 | 6 |
| MultipartDef | 0 | 0 | 0 | 3 | 3 |
| NamingResources | 2 | 0 | 0 | 0 | 2 |
| SecurityCollection | 2 | 0 | 0 | 0 | 2 |
| ServletDef | 4 | 0 | 0 | 0 | 4 |
| CsrPreventionFilter | 4 | 0 | 0 | 0 | 4 |
| AbstractReplicatedMap | 2 | 0 | 0 | 0 | 2 |
| FeatureInfo | 2 | 1 | 0 | 0 | 3 |
| DeltaSession | 0 | 1 | 0 | 0 | 1 |
| All | 19 | 4 | 2 | 6 | 31 |

TABLE 6
SMOs Found across Five Versions of Eiffel

| | attribute added | attribute removed | attribute type changed | inheritance added | inheritance removed | class name changed | void safety added | all |
|-----------|-----------------|-------------------|------------------------|-------------------|---------------------|--------------------|-------------------|-----|
| 1993-2002 | 8 | 2 | 0 | 8 | 14 | 0 | 0 | 32 |
| 2002-2006 | 2 | 0 | 3 | 4 | 0 | 0 | 0 | 9 |
| 2006-2008 | 4 | 1 | 0 | 3 | 0 | 1 | 0 | 9 |
| 2008-2010 | 0 | 0 | 0 | 0 | 0 | 0 | 4 | 4 |
| All | 14 | 3 | 3 | 15 | 14 | 1 | 4 | 54 |

time period on which we focus is wider, ranging from 1993 to 2010 instead of from 1999 to 2009. The main reason is that given that Eiffel is less widespread than Java, we hope to find more interesting data in a wider time interval. In both cases we allow both languages to have an initial stabilization period. The availability of data in the software repositories is also a limiting factor. The list of relevant code changes for Eiffel persistence allows us to extend the previous list by adding "inheritance from one class added," "inheritance from one class removed," "class name changed," and "Void safety clause added." With respect to the Java study, we do not consider "Attribute initialization value changed" because in Eiffel assigning a value to an attribute in the declaration means making it a constant, which is typically not serialized. Furthermore, we have not found any instance of "Attribute renamed," "Attribute becoming a constant," or "Constant becoming an attribute." Interestingly, however, we find a consistent number of inheritance-related code changes, of which there was no trace in Java. This may be due to the fact that the Eiffel object model allows multiple inheritance, making it easier to add and remove functionality by using inheritance itself. As we model inheritance by considering the flattened class schema, the "Inheritance added" and "Inheritance removed" code changes fall back into the other already analyzed ones. Finally, there are also some changes to attributes meant to enforce void-safety, a matter discussed in a separate paragraph below. The kinds of code changes found, together with their occurrences and across considered versions, are synthesized in Table 6. Note that "Attribute added" and "Attribute removed" are still relevant, while "Inheritance added" and "Inheritance removed" together constitute 54 percent of all persistence-related code changes.

Table 7 shows the distribution of code changes across different classes and across all the considered five versions of Eiffel. As in the Java study, the data suggest that persistence-related code changes are sufficiently widespread among classes.

During our analysis we found a particularly interesting example of the kind of issues ESCHER would help solving. In class `HASH_TABLE`, considered from revision 13,752 to revision 47,039, there was a switch to a more efficient container by changing the types of three attributes. This resulted in the previously serialized objects not being readable anymore. The solution, implemented in a later version, was to code a transformation function using an ad

TABLE 7
SMOs Found across Five Eiffel Versions, per Class

| Class | Persistence-related changes |
|------------------------|-----------------------------|
| Arrayed_list | 2 |
| Bit_ref | 0 |
| Date_time_parser | 6 |
| Array | 1 |
| Date | 2 |
| String_8 | 6 |
| Tuple | 3 |
| Action_sequence | 2 |
| Array_2 | 0 |
| Arrayed_set | 0 |
| Hash_table | 12 |
| Linked_tree | 1 |
| Linked_set | 2 |
| Linked_list | 0 |
| I18n_locale_manager | 4 |
| Priority_queue | 0 |
| Random | 3 |
| Time | 2 |
| Binary_search_tree | 3 |
| Binary_search_tree_set | 4 |
| Uuid | 1 |
| Arrayed_stack | 0 |
| All | 54 |

hoc created attribute, *hash_table_version_57*. If the library maintainers had used ESCHER, there would have been no need to add any attribute, and the transformation function would not have been cluttering the class.

Void safety-related code changes. When we refer to void-safety we mean the guarantee that, in an object-oriented program, the operation $x.f(\dots)$ will never fail because of the target x denoting a void reference at execution time [42]. The Eiffel language incorporates the mechanism into the compiler, requiring developers to choose, for every attribute, one keyword as a default (*attached* or *detachable*) and to specify the other when needed. In the case of the class *DATE_TIME_PARSER*, we have discovered four code changes related to void safety in which the keyword *detachable* was added to four attributes. This code change (from *attached* to *detachable*) is harmless under the point of view of software evolution because an attached attribute, which is forced to be nonvoid, will always be stored as such and so it will comply with the requirement of being detachable at retrieval time. Though it did not occur in our findings, the symmetric code change, from *detachable* to *attached*, can in principle be dangerous. A detachable attribute stored as void will be considered as attached by the retrieving system. This will raise an exception at runtime, both in case of retrieval from storage and in any other case involving assignment at runtime, such as a call to an external module assigning void to an attached attribute. We believe that this finding, while related to the Eiffel language, can still be useful for future evolution of languages like Java that have not yet incorporated this feature.

5.2.2 Analyzing the EiffelStudio Code Base

We choose to analyze EiffelStudio as an example of an IDE, a very common type of application that is likely to have persistence needs. EiffelStudio stores information about user preferences and also stores data concerning automatically generated tests. The EiffelStudio IDE was initially called EiffelBench and was released in 1990. It was renamed

TABLE 8
SMOs Found in EiffelStudio's Repository

| | attribute added | attribute removed | attribute renamed | attribute type changed | inheritance added | class name changed | void safety added | all |
|-------------------------|-----------------|-------------------|-------------------|------------------------|-------------------|--------------------|-------------------|-----|
| Sd_config_data | 8 | 0 | 1 | 0 | 0 | 0 | 0 | 9 |
| Search_table | 2 | 0 | 0 | 0 | 0 | 0 | 1 | 3 |
| Sd_inner_container_data | 0 | 0 | 0 | 2 | 0 | 0 | 4 | 6 |
| Ev_split_area | 1 | 0 | 0 | 0 | 0 | 0 | 2 | 3 |
| Ev_split_area_i | 0 | 1 | 0 | 0 | 0 | 0 | 3 | 4 |
| Ev_widget_i | 1 | 0 | 0 | 0 | 0 | 0 | 3 | 4 |
| Sd_place_holder_zone | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| Sd_notebook_upper | 4 | 3 | 0 | 2 | 0 | 0 | 0 | 9 |
| Epa_test_case_info | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 |
| Profile_information | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 |
| Profile_set | 0 | 0 | 0 | 36 | 0 | 0 | 0 | 36 |
| Test_session_record | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 |
| All | 17 | 4 | 1 | 41 | 1 | 1 | 13 | 78 |

to EiffelStudio in 2001 and was open sourced in 2006. We analyze the evolution of persistent classes of the main EiffelStudio application code and of its automated testing facility called "AutoTest." We take into account the earliest and the latest releases available in the open source repository at the time of writing (5.4 from 2004 and 6.8 from 2011).⁷ In Table 8 we show, for brevity, the classes in which we actually found at least one SMO. This occurred in 12 out of the 20 classes analyzed. The data once again confirm what we found in the case of the Eiffel library classes, and are also consistent with what we found in the Java repositories.

6 IMPLEMENTATION

The ESCHER tool implements the model presented in Section 4 making it usable in practice. It is integrated into EiffelStudio, and can also be used independently as a library. While our approach can be applied to any object-oriented programming language providing support for storing and retrieving objects, we use Eiffel (and its IDE EiffelStudio) because of its integrated support for Design by Contract, and in particular for class invariants [43]. The detailed documentation about ESCHER, a step-by-step illustrated tutorial, the source code, and the executables for the tool are available for download on the ESCHER project page.⁸

6.1 Version Handling

ESCHER handles class versions either by using the existing version information or by automatically assigning a version tag if not present. The versioning mechanism relies on the notion of *release*: a versioned, compiled, and semantically coherent set of classes constituting a software system. Different versions of the same class can only be part of different releases. At release time, ESCHER automatically updates the class version tags of the classes that actually changed. The release version is increased if at least one class has changed with respect to the previous release. Provided

7. <http://svn.eiffel.com/eiffelstudio/trunk/> (last visited: 15/6/12).

8. <http://trac.inf.ethz.ch/trac/meyer/eve/wiki> (last visited: 15/6/12).

the required transformation functions exist, it will always be possible to attempt a retrieval of any object of a certain class and version into an object of another version of the same class. This applies to both forward and backward transformations as customers might be running an old system (specified by an old release number) and in need of retrieving objects stored by a newer system release.

6.2 Code Generation

While modest, the support for generating code provided by ESCHER relieves developers from writing boilerplate code to transform instances stored in different versions with respect to the current one. It also lets developers focus on the actual specification of the transformation function rather than on the framework code necessary to wire all the transformation functions. Given a class *C*, ESCHER instruments it with a version tag, and automatically creates a corresponding handler class containing all the related transformation functions. The generated name of the handler class is `C_SCHEMA_EVOLUTION_HANDLER` in this case. All the associations between each handler name and its corresponding class are kept in the ad hoc generated class `SCHEMA_EVOLUTION_PROJECT_MANAGER`. ESCHER automatically creates and manages all the separate folders containing the various releases as well. The code generation algorithm for the transformation functions compares the two abstract syntax trees of two class versions statically, looking for known SMOs. If a known SMO is detected, the algorithm generates a transformation function and code according to the respective heuristic rule. The transformation function is then placed into an ad hoc handler class, containing all the existing transformation functions between pairs of different versions. This helps to avoid polluting the class code with evolution-related code, and makes the whole approach more scalable. Comments, hints, and informational messages are generated to guide developers through the process.

6.3 IDE Support and Runtime Mechanism

ESCHER is seamlessly integrated into the EiffelStudio IDE, and its functionalities are accessible through an ad hoc panel. In addition to release a project, developers can trigger the creation and release of a transformation function by selecting two versions from a dynamically generated panel that only lists the versions actually existing in the repository. Developers can also create a “filter” that allows them to select the attributes they want to be part of the serialized form.

At retrieval time, ESCHER relies on an ad hoc algorithm automatically performing conversions across two different versions. The algorithm will raise an exception either when a specific schema evolution handler does not exist, or when a specific transformation function between the two versions does not exist. The class invariant check for the retrieving object is also enforced.

7 ASSESSMENT AND FUTURE WORK

7.1 Evaluating ESCHER on the Java and Eiffel Repositories

Analyzing the history of classes along different versions provides useful insights on the schema evolution process. The data analysis presented in Section 5 shows that the

schema of persistent classes may evolve significantly both in software libraries and in applications. To evaluate ESCHER, we check if it recognizes the 278 persistence-related SMOs (counting both Eiffel and Java SMOs) found in our analysis. The results are that ESCHER is able to identify the SMOs in 93 percent of all cases (259 out of 278). The unrecognized SMOs are “Void safety added” and “Class name changed,” which is interpreted as a new class introduced into the system. We count “Attribute renamed” as recognized because ESCHER, though not generating the full body of the transformation function, is aware of the different possibilities and generates a comment providing guidance to developers. We also observe that in none of the four code repositories analyzed was an “Attribute renamed” a false positive, that is, had to be interpreted as an “Attribute removed” plus “Attribute added.” The results above suggest that the simple model we devised is sufficiently detailed to take into account a realistic set of SMOs in two different object-oriented programming languages.

We have shown in Section 4 that the transformations functions are complete with respect to attribute modifications. With respect to their soundness, we do not provide static guarantees, but given any two versions and given a transformation function between the two, the tool ensures that the function preconditions and postconditions, and the retrieving class invariant, hold at runtime.

We have already remarked on the importance of class invariants to avoid having potentially unsafe objects accepted into the system. Unfortunately, class invariants are not yet widely accepted among developers, and not all mainstream languages offer native support for them. An obvious question is then: What can Java developers do to emulate the ESCHER mechanism enforcing class invariants? While analyzing the Java libraries, we discovered a possible solution in class *BitSet* during the transition from version 5.0 to 6.0. A new method `checkInvariants()` is introduced, with the idea of enforcing some class-wide properties. This method is then invoked from every public method in the class, emulating the Eiffel invariant checking mechanism. In conclusion, the absence of an embedded Design By Contract language support (also for preconditions and postconditions) makes things more complicated for Java developers, but the example of class *BitSet* shows that emulating the invariant mechanism remains possible.

In the worst case, developers can still bypass, actively or passively, all the checks mentioned. In case of an attribute added, for example, they could ignore the warning, accept the (possibly wrong) default provided, and in general fail to code an appropriate transformation function or class invariant. They could even code the right invariant but disable runtime checks for assertions. This means that while ESCHER supports developers in handling the class schema evolution of an object-oriented software application, it does not guarantee a completely safe class schema evolution. To improve this, we believe that our measure for evolution robustness can help in assessing how much trust we can put in a code release with respect to class schema evolution.

7.2 Threats to Validity

While the analysis of the results of our four studies shows that the ESCHER model and its implementation address most practical challenges, the validation based on empirical

data is by definition questionable. In addition, the current SMO detection and code generation could be dependent on the specific programming languages chosen (Java and Eiffel).

Assuming that the idea and the current implementation are meaningful, general willingness to adopt ESCHER remains to be proven. In particular, there is no usability study trying to investigate developers' understanding and acceptance of such a solution. We did not perform a usability study because we gave priority to validating ESCHER using existing, realistic, and widely known software projects.

7.3 Conclusions

This work proposes a shift of attitude with respect to how developers cope with class schema evolution. By integrating support into an IDE, the approach elevates class schema evolution to the status of first-class citizen of the software development process rather than undesirable side effect of the software production activities. It also proposes a significant time shift (from runtime to release time) for the evolution-related activities. A system release then becomes an event triggering a whole set of tool actions intended to help developers to better focus on possible issues that may arise from previous versions of the newly released code.

Possible future work is to broaden the empirical validation by performing studies on more code repositories, in particular in object-oriented programming languages other than Java and Eiffel. We also plan to integrate object browsing and writing-back capabilities into the IDE to help developers having to analyze objects previously serialized using a non-human-readable format. Finally we plan to consider, in our model of software evolution, the evolution of the persistence software mechanisms themselves.

ACKNOWLEDGMENTS

The authors thank the Swiss National Foundation (grant numbers 200021-113808) for partially funding this work. They also thank Matthias Loeu and Teseo Schneider for their implementation work on the initial versions of ESCHER.

REFERENCES

- [1] M. Piccioni, M. Oriol, B. Meyer, and T. Schneider, "An Ide-Based, Integrated Solution to Schema Evolution of Object-Oriented Software," *Proc. IEEE/ACM Int'l Conf. Automated Software Eng.*, pp. 650-654, 2009.
- [2] M.P. Atkinson, "Persistence and Java—A Balancing Act," *Proc. Int'l Symp. Objects and Databases*, pp. 1-31, 2000.
- [3] M. Jordan, "A Comparative Study of Persistence Mechanisms for the Java Platform," Technical Report TR-2004-136, Sun Microsystems Laboratories, 2004.
- [4] Oracle, "Java Serialization Specification," <http://download.oracle.com/javase/6/docs/platform/serialization/spec/version.html#6519>, 2012.
- [5] Microsoft, "Version Tolerant Serialization," [http://msdn.microsoft.com/en-us/library/ms229752\(VS.100\).aspx](http://msdn.microsoft.com/en-us/library/ms229752(VS.100).aspx), 2012.
- [6] B. Meyer, "Applying 'Design by Contract'," *Computer*, vol. 25, no. 10, pp. 40-51, 1992.
- [7] J. Paterson, S. Edlich, H. Hörning, and R. Hörning, *Microsoft Windows Internals*. Microsoft Press, 2005.
- [8] J. Banerjee, W. Kim, H.-J. Kim, and H.F. Korth, "Semantics and Implementation of Schema Evolution in Object-Oriented Databases," *Proc. ACM SIGMOD Int'l Conf. Management of Data*, pp. 311-322, 1987.
- [9] A.H. Skarra and S.B. Zdonik, "The Management of Changing Types in an Object-Oriented Database," *Proc. Conf. Object-Oriented Programming Systems, Languages and Applications*, pp. 483-495, 1986.
- [10] S.-E. Lautemann, "An Introduction to Schema Versioning in Oodbms," *Proc. Seventh Int'l Workshop Database and Expert Systems Applications*, pp. 132-139, 1996.
- [11] S.R. Monk and I. Sommerville, "Schema Evolution in OODBs Using Class Versioning," *SIGMOD Record*, vol. 22, no. 3, pp. 16-22, 1993.
- [12] GemStone, "Gemstone/s 64 Bit Programming Guide," <http://www.gemstone.com/docs/GemStoneS/GemStone64Bit/2.4.4.3/GS64-ProgGuide-2.4.pdf>, 2009.
- [13] Oracle "Edition-Based Redefinition," <http://www.oracle.com/technetwork/database/features/availability/ed.-based-redefinition-1-133045.pdf>, 2009.
- [14] M. Piccioni, M. Oriol, and B. Meyer, "Ide-Integrated Support for Schema Evolution in Object-Oriented Applications," *Proc. Workshop Reflection AOP and MetaData for Software Evolution*, pp. 27-36, 2007.
- [15] B. Meyer, "Conversions in an Object-Oriented Language with Inheritance," *J. Object-Oriented Programming*, vol. 13, no. 9, pp. 28-31, 2001.
- [16] G.M. Bierman, M.J. Parkinson, and J. Noble, "Upgradej: Incremental Typechecking for Class Upgrades," *Proc. 22nd European Conf. Object-Oriented Programming*, pp. 235-259, 2008.
- [17] D. Dig, K. Manzoor, R. Johnson, and T.N. Nguyen, "Refactoring-Aware Configuration Management for Object-Oriented Programs," *Proc. 29th Int'l Conf. Software Eng.*, pp. 427-436, 2007.
- [18] K. Czarnecki, J.N. Foster, Z. Hu, R. Lämmel, A. Schürr, and J.F. Terwilliger, "Bidirectional Transformations: A Cross-Discipline Perspective," *Proc. Second Int'l Conf. Theory and Practice of Model Transformations*, pp. 260-283, 2009.
- [19] M. Hofmann, B.C. Pierce, and D. Wagner, "Symmetric Lenses," *Proc. ACM SIGPLAN-SIGACT Symp. Principles of Programming Languages*, 2011.
- [20] M.P. Atkinson, "Programming Languages and Databases," *Proc. Fourth Int'l Conf. Very Large Data*, pp. 408-419, 1978.
- [21] M.P. Atkinson and R. Morrison, "Orthogonally Persistent Object Systems," *Int'l J. Very Large Data Bases*, vol. 4, no. 3, pp. 319-401, 1995.
- [22] M.P. Atkinson and M. Jordan, "A Review of the Rationale and Architectures of Pjama: A Durable, Flexible, Evolvable and Scalable Orthogonally Persistent Programming Platform," Technical Report SMLI TR-2000-90, Sun Microsystems Laboratories, 2000.
- [23] B.M.B. Lewis and N. Gafter, "Architecture of the PEVM: A High-Performance Orthogonally Persistent Java Virtual Machine," Technical Report TR-2000-93, Sun Microsystems Laboratories, 2000.
- [24] M. Dmitriev, "Safe Class and Data Evolution in Long-Lived Java Applications," Technical Report SMLI TR-2001-98, Sun Microsystems Laboratories, 2001.
- [25] G. Wachsmuth, "Metamodel Adaptation and Model Co-Adaptation," *Proc. 21st European Conf. Object-Oriented Programming*, pp. 600-624, 2007.
- [26] B. Gruschko, D.S. Kolovos, and R.F. Paige, "Towards Synchronizing Models with Evolving Metamodels," *Proc. Int'l Workshop Model-Driven Software Evolution*, 2007.
- [27] M. Herrmannsdoerfer, S. Benz, and E. Jürgens, "Automatability of Coupled Evolution of Metamodels and Models in Practice," *Proc. 11th Int'l Conf. Model Driven Eng. Languages and Systems*, pp. 645-659, 2008.
- [28] A. Cicchetti, D.D. Ruscio, R. Eramo, and A. Pierantonio, "Automating Co-Evolution in Model-Driven Engineering," *Proc. 12th IEEE Int'l Enterprise Distributed Object Computing Conf.*, pp. 222-231, 2008.
- [29] J.-R. Falleri, M. Huchard, M. Lafourcade, and C. Nebut, "Metamodel Matching for Automatic Model Transformation Generation," *Proc. 11th Int'l Conf. Model Driven Eng. Languages and Systems*, pp. 326-340, 2008.
- [30] I. Neamtii, M.W. Hicks, G. Stoye, and M. Oriol, "Practical Dynamic Software Updating for C," *Proc. ACM SIGPLAN Conf. Programming Language Design and Implementation*, pp. 72-83, 2006.
- [31] R.G.G. Cattell, *The Object Data Standard: ODMG 3.0*. Morgan Kaufmann, 2000.

- [32] C. Delgado, J. Samos, and M. Torres, "Primitive Operations for Schema Evolution in ODMG Databases," *Proc. Object-Oriented Information Systems*, pp. 226-237, 2003.
- [33] C. Date, *Introduction to Database Systems*, eighth ed. Addison Wesley, 2003.
- [34] E. Rahm and P.A. Bernstein, "A Survey of Approaches to Automatic Schema Matching," *Int'l J. Very Large Data Bases*, vol. 10, no. 4, pp. 334-350, 2001.
- [35] T. Milo and S. Zohar, "Using Schema Matching to Simplify Heterogeneous Data Translation," *Proc. 24th Int'l Conf. Very Large Data Bases*, pp. 122-133, 1998.
- [36] B.S. Lerner, "A Model for Compound Type Changes Encountered in Schema Evolution," *ACM Trans. Database Systems*, vol. 25, no. 1, pp. 83-127, 2000.
- [37] C.A. Curino, H.J. Moon, and C. Zaniolo, "Graceful Database Schema Evolution: The Prism Workbench," *Proc. VLDB Endowment*, vol. 1, pp. 761-772, 2008.
- [38] V. Authors, "Schema Evolution Benchmark," <http://schemaevolution.org/>, 2012.
- [39] C. Boyapati, B. Liskov, L. Shriram, C.-H. Moh, and S. Richman, "Lazy Modular Upgrades in Persistent Object Stores," *Proc. 18th Ann. ACM SIGPLAN Conf. Object-Oriented Programming, Systems, Languages, and Applications*, pp. 403-417, 2003.
- [40] E. Franconi, F. Grandi, and F. Mandreoli, "Schema Evolution and Versioning: A Logical and Computational Characterisation," *Proc. Workshop Foundations of Models and Languages for Data and Objects*, pp. 85-99, 2000.
- [41] D. Advani, Y. Hassoun, and S. Counsell, "Extracting Refactoring Trends from Open-Source Software and a Possible Solution to the "Related Refactoring" Conundrum," *Proc. ACM Symp. Applied Computing*, pp. 1713-1720, 2006.
- [42] B. Meyer, "Attached Types and Their Application to Three Open Problems of Object-Oriented Programming," *Proc. 19th European Conf. Object-Oriented Programming*, pp. 1-32, 2005.
- [43] *ECMA International Standard: Eiffel Analysis, Design and Programming Language*, ECMA Committee TC39-TG4 Std. 367, 2005.



Marco Piccioni received the degree in mathematics from the Università degli Studi di Roma La Sapienza, and the master's degree in economics from the Università L. Bocconi, Milano. After this, he worked for 10 years at Sistemi Informativi S.p.A. (an IBM company) as a software engineer and trainer. He is a research assistant and also working toward the PhD degree at the Chair of Software Engineering at ETH Zurich. His research work at ETH focuses on the software evolution of persistent object-oriented applications. He is a member of the IEEE.



Manuel Oriol graduated from ENSEEIHT school with an engineer diploma in computer science and received the MSc degree in artificial intelligence from INP Toulouse. He subsequently graduated from the University of Geneva with the PhD degree in information systems. He is working as a principal scientist with ABB Corporate Research, Industrial Software Systems in Baden-Dättwil, Switzerland. He is also a lecturer at the University of York in the United Kingdom. His research interests include a large spectrum and include object-oriented languages, software testing, component infrastructure, computer science education, middleware/cloud, real-time systems, and dynamic software updating.



Bertrand Meyer is working as a professor of software engineering at ETH Zurich, as head of the Software Engineering Laboratory at ITMO, Saint Petersburg, Russia, and as chief architect of Eiffel Software.

► For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/publications/dlib.