

Do software models based on the UML aid in source-code comprehensibility? Aggregating evidence from 12 controlled experiments

Giuseppe Scanniello¹ · Carmine Gravino² ·
Marcela Genero³ · José A. Cruz-Lemus³ ·
Genoveffa Tortora² · Michele Risi² · Gabriella Doderò⁴

© Springer Science+Business Media, LLC, part of Springer Nature 2018

Abstract In this paper, we present the results of long-term research conducted in order to study the contribution made by software models based on the Unified Modeling Language (UML) to the comprehensibility of Java source-code deprived of comments. We have conducted 12 controlled experiments in different experimental contexts and on different sites with participants with different levels of expertise (i.e., Bachelor's, Master's, and PhD students and software practitioners from Italy and Spain). A total of 333 observations were obtained from these experiments. The UML models in our experiments were those produced

Communicated by: Richard Paige

✉ Carmine Gravino
gravino@unisa.it

Giuseppe Scanniello
giuseppe.scanniello@unibas.it

Marcela Genero
marcela.genero@uclm.es

José A. Cruz-Lemus
joseantonio.cruz@uclm.es

Genoveffa Tortora
tortora@unisa.it

Michele Risi
mrisi@unisa.it

Gabriella Doderò
Gabriella.Doderò@unibz.it

¹ University of Basilicata, Potenza, Italy

² University of Salerno, Fisciano, Italy

³ University of Castilla-La Mancha, Ciudad Real, Spain

⁴ Free University of Bozen, Bolzano-Bozen, Italy

in the analysis and design phases. The models produced in the analysis phase were created with the objective of abstracting the environment in which the software will work (i.e., the problem domain), while those produced in the design phase were created with the goal of abstracting implementation aspects of the software (i.e., the solution/application domain). Source-code comprehensibility was assessed with regard to correctness of understanding, time taken to accomplish the comprehension tasks, and efficiency as regards accomplishing those tasks. In order to study the global effect of UML models on source-code comprehensibility, we aggregated results from the individual experiments using a meta-analysis. We made every effort to account for the heterogeneity of our experiments when aggregating the results obtained from them. The overall results suggest that the use of UML models affects the comprehensibility of source-code, when it is deprived of comments. Indeed, models produced in the analysis phase might reduce source-code comprehensibility, while increasing the time taken to complete comprehension tasks. That is, browsing source code and this kind of models together negatively impacts on the time taken to complete comprehension tasks without having a positive effect on the comprehensibility of source code. One plausible justification for this is that the UML models produced in the analysis phase focus on the problem domain. That is, models produced in the analysis phase say nothing about source code and there should be no expectation that they would, in any way, be beneficial to comprehensibility. On the other hand, UML models produced in the design phase improve source-code comprehensibility. One possible justification for this result is that models produced in the design phase are more focused on implementation details. Therefore, although the participants had more material to read and browse, this additional effort was paid back in the form of an improved comprehension of source code.

Keywords Aggregation · Heterogeneity · Unified modeling language · Controlled experiments

1 Introduction

The Unified Modeling Language (UML) (OMG 2005) is considered to be the de-facto standard in the analysis, design, and evolution of object-oriented software (Erickson and Siau 2007; Grossman et al. 2005), despite the fact that domain-specific *modeling* languages are increasing in popularity (Hutchinson et al. 2011). However, many software companies are still reluctant to use UML because it is perceived to be difficult to learn and use Agarwal and Sinha (2003). It may, therefore, be important to investigate whether or not the use of the UML makes a practical difference in software development and evolution, thus possibly encouraging resilient companies to adopt UML.

The UML has been the subject of a number of empirical studies in the software engineering field (Budgen et al. 2011). Of these studies, only a few are focused on the usage of this notation throughout the software development life cycle (Anda et al. 2006). This lack is even more evident in software maintenance and evolution (e.g., Bavota et al. 2013, Scanniello et al. 2014). Scanniello et al. (2010) conducted an industrial survey regarding the use of the UML in the software industry. The results showed that software engineers wishing to deal with software maintenance and evolution very often have at their disposal only UML-based models (or simply UML models) produced in the requirements engineering process (or analysis phase) and, in a few cases, those produced in the design phase. Using the findings of this survey as a basis, we began long-term research with the aim of studying the contribution made by UML-based models produced in the analysis and design phases to source-code comprehensibility. In particular, we conducted 12 controlled experiments with

different kinds of participants (Gravino et al. 2010, 2015; Scanniello et al. 2010, 2014, 2015) to study the effect of these kinds of models on the comprehension of source code deprived of comments. We removed comments to avoid their effect being confused with the main factor studied (i.e., the presence or the absence of software models).

In order to aggregate the results of these experiments and to obtain the global effect of analysis and design UML models on source-code comprehensibility, we carried out a meta-analysis on the results obtained from the individual experiments. In this paper, we present the results of this aggregation by attempting to answer the following research question:

- *Do software models produced in the analysis and design phases aid the comprehension of Java source code (assessed on the basis of the answers developers provide to questions on that code), and do these models affect the time taken to comprehend that code and the efficiency with which that comprehension occurs?*

The research work presented in this paper is based on that presented in Scanniello et al. (2015), in which we showed the preliminary results obtained after the aggregation of our 12 controlled experiments. The current paper extends the previous work as follows: (i) we have considered issues related to the heterogeneity of the individual experiments in the aggregation of their results; (ii) we have considered an additional dependent variable, namely Efficiency, which is computed as the ratio between the level of comprehension achieved when performing a task and the time required to complete it; and (iii) we have extended the discussion concerning related work, experimental results, and threats to validity.

In this new paper, we provide a detailed description of the following main contributions:

- The results of global effect of analysis and design UML models on source-code comprehensibility;
- A discussion regarding the possible practical implications of the results of our study;
- How to deal with the heterogeneity of experiments when using a meta-analysis to aggregate the results obtained from them.

This paper is organized as follows. We discuss related work in Section 2, while the background is presented in Section 3. Our long-term research is presented in Section 4, while the results obtained are shown and discussed in Section 5. In Section 5, we also discuss the practical implications of the results of our study from the perspectives of both researchers and professionals. We conclude the paper with our final remarks and future work in Section 6.

2 Related work

In accordance with the research question stated previously, in this section we first focus on a set of works which highlights the use of UML diagrams as a means to maintain source-code. We conclude this section by presenting research work on model-based traceability, since the use of traceability links might be a viable means to support the comprehension and maintenance of source code.

2.1 UML and software maintenance

There are two literature reviews related to the topic that should be taken into account (Fernández-Sáez et al. 2013; Zhi et al. 2015). First, Fernández-Sáez et al. (2013) present a systematic mapping study (SMS) whose goal was to discover the empirical evidence related to the use of UML diagrams in source-code maintenance and the maintenance of

UML diagrams themselves. A total of 38 papers were found as a result of this SMS, including 66 empirical studies, but only two of them were specifically focused on source-code maintenance (Arisholm et al. 2006; Dzidek et al. 2008). Zhi et al. (2015) also report another SMS but with a slightly different goal: the existing literature concerning software documentation cost, benefit and quality. Once again, the same two works (Arisholm et al. 2006; Dzidek et al. 2008) are directly related to source-code maintenance. These two empirical studies will be explained below, together with another set of experiments which was published after the performance of the two SMSs, and thus not included in them.

Using the aforementioned work as a basis, Dzidek et al. (2008) performed an experiment with a group of 20 professionals. Half of them used UML documentation to carry out a set of modifications on a web-based system developed in Java, while the other half did not. This was done to assess whether providing UML documentation reduced the effort required to correctly implement a set of tasks regarding changes and increased the functional correctness and design quality of those changes. The results reported indicate that UML users had to spend more time, especially when the documentation had to be updated. However, the use of UML documentation was simultaneously always beneficial in terms of functional correctness, since fewer faults were introduced into the software maintained. It is also important to highlight that the no-UML group of participants had more problems as regards understanding the most complex part of the system.

Arisholm et al. (2006) conducted two controlled experiments to assess whether UML documentation helped to reduce the effort required to change the source-code of a software system. The original experiment was conducted in Oslo, Norway, and the second in Ottawa, Canada. The first experiment was conducted with 20 3rd year undergraduate students, while the second was carried out with 4th year 78 undergraduate students. In both these experiments, the independent variable was: using (or not using) UML documentation. The participants' performances were measured by considering the time required to perform changes, excluding and including diagram modifications, correctness of these changes and the quality of the changed design. The most important result obtained was: UML documentation does not provide an advantage as regards time, although it helps to improve the correctness and quality achieved when solving the most complex tasks.

Leotta et al. (2013) present a pilot experiment carried out to relate the level of alignment between UML documents and code and maintainers' efficiency. A group of 21 undergraduate students had to perform a set of 4 maintenance tasks on two systems using the Eclipse framework while surfing the UML documents provided (sequence and class diagrams). Although this was a pilot study, the results confirmed the general belief that a more aligned documentation is of greater assistance during maintenance tasks.

Fernández-Sáez et al. (2015) report a family of experiments carried out at two different universities in Italy and Spain. The aim of this study was to assess whether the origin of UML diagrams (the design phase of a life-cycle or a reverse engineering technique) influences the maintenance of the corresponding source-code. A controlled experiment and two replications of it were performed. The authors involved a total of 149 MSc students (categorized according to their ability, which was calculated using their course grades). These students had to carry out adaptive and corrective maintenance tasks. The main finding of the work was that participants with a higher ability achieved better scores when using the diagrams with a forward design origin, while low ability participants got better scores when using reverse engineered diagrams. The authors provide a possible explanation for this situation by relating low levels of experience to difficulty in using (and/or understanding) UML diagrams.

The last experiment in this list Fernández-Sáez et al. (2016) reports a family of four controlled experiments carried out by over 80 BSc and MSc students at three different universities in Italy, The Netherlands, and Spain. In this case, the goal was to analyze whether

a high or a low level of detail in UML diagrams has an impact on the maintainability of source-code during a model-centric development. The maintainability of the source-code was measured by means of its understandability and modifiability, and the participants had to answer a multiple choice questionnaire to show how they had understood the system, in addition to performing a set of corrective maintenance tasks to show how it could be modified. The results of the family of experiments indicated that diagrams with a high level of details improved the understandability of the system, while those with a low level of details improved its modifiability. However, the authors indicate that, when attempting to replicate the study, the results obtained seemed to be incoherent and to have no clear tendency. After a thorough evaluation of this, they discovered that the diagrams had possibly been used incorrectly, or even not used at all, which would have led to the results obtained. They relate this to the day-by-day situation in industry, when users' false self-certainty may lead them not to use (or not to use properly) the system documentation, and eventually find themselves involved in unexpected and undesired situations as regards their maintenance duty.

Table 1 presents a summary of the main features of these papers. The columns in the table are labeled as:

- Ref: contains the bibliographical reference to the paper.
- Goal: describes the objective of the experiment.
- Participants: presents the number of participants that took part in the experiment/s, plus their type (students, professionals, etc.).
- Independent variable: describes the variable whose effect on the dependent variables is to be studied. The values (treatments) of the independent variable are also presented.
- Dependent variables: presents the outcome variables, i.e., those affected by the changes made to the independent variables.
- Tasks: describes the tasks to be performed by the participants as part of the experiment.
- Results: shows the main findings obtained in the experiment.

Some other empirical studies in industry related to the research question of this work can also be found. For example, Scanniello et al. (2010) show the results of an explorative survey carried out to investigate the state of practice as regards the use of UML in software development and maintenance. At that time, UML appeared to be the most widely used modeling option for software development and maintenance. In particular, 74% and 75% of the companies interviewed employed UML in the development and maintenance phases, respectively. The authors did not ask the remaining survey respondents which modeling languages were used as an alternative to UML. The companies interviewed stated that maintenance operations are commonly performed by practitioners who do not have a vast amount of experience. According to the study, an ordinary maintenance operation, such as making certain corrective changes, needs an effort range of between 1 and 5 persons/hour. These values are multiplied by 10 in the case of an extraordinary maintenance operation, such as perfective or adaptive changes.

Fernández-Sáez et al. (2013) also present the preliminary findings of a case study whose intention was to discover whether the investment in UML is justified by the benefits (improved productivity and product quality) in software maintenance projects. They consequently focus on discovering what the cost and the payback of using UML in a software maintenance project are. They carried out a case study in a multinational company with an IT department of 800-1000 employees. They collected data from the files shared by the department and, principally, by interviewing the company's personnel (20 useful interviews). They eventually concluded that the employees reported several benefits of using UML: a better understanding of the problem domain, improved communication, a reduction in

Table 1 Summary of experiments features

Ref.	Goal	Participants	Independent variable	Dependent variables	Tasks	Results
Arisholm et al. (2006)	Assessing whether the UML documentation provided helped to reduce the effort required to change the source-code of a software system.	20 + 78 undergraduate (3rd and 4th year) students	Using (or not) UML documentation	Time taken to perform changes, excluding including diagram modifications, correctness of these changes and quality of the changed design.	Performing several operations related to the systems used.	UML documentation does not provide an advantage as regards time, although it helps to improve the correctness and quality achieved when solving the most complex tasks.
Dzidek et al. (2008)	Assessing whether providing UML documentation reduced the effort required to correctly implement a set of change tasks, and increased the functional correctness and design quality of those changes.	20 professional developers	Using (or not) UML documentation	Time taken to perform changes, excluding including diagram modifications, correctness of these changes (in terms of the number of submissions of a solution with a fault and their relation to the existing functionalities) and quality of the changed design.	Modifying the existing classes or adding new ones to change the system's behavior (adding new functionalities, changing the original ones, etc.).	UML participants spent more time than the others updating documentation, but using UML was beneficial in terms of functional correctness and helped participants to understand the systems better (especially the more complex ones).

Table 1 (continued)

Leotta et al. (2013)	<p>Relating the level of alignment between UML documents and code and maintainers' efficiency.</p>	<p>21 undergraduate students</p>	<p>Document-tation alignment (values more and less)</p>	<p>Efficiency as regards performing maintenance tasks.</p>	<p>Execution of four maintenance tasks on a system.</p>	<p>A more aligned documentation is more helpful during maintenance tasks.</p>
Fernández-Sáez et al. (2015)	<p>Investigating whether the origin of UML diagrams influences the maintenance of the corresponding source-code.</p>	<p>20+51+78 MSC students</p>	<p>Origin of the diagrams (forward designed or reverse engineered)</p>	<p>Source-code maintainability, measured by its effectiveness and efficiency.</p>	<p>Five different tasks including adaptive and corrective maintenance tasks.</p>	<p>Participants with a higher ability achieved better scores using the diagrams with a forward design origin, possibly because UML helps more experienced users.</p>
Fernández-Sáez et al. (2016)	<p>Analyzing the level of detail in UML diagrams and checking its impact on the maintainability of source-code during a model-centric development.</p>	<p>65 BSc + 16 MSC students</p>	<p>Level of detail of the UML diagrams (high or low)</p>	<p>Modifiability and understanding of the source-code, both measured by their effectiveness and efficiency.</p>	<p>Answering 3 multiple choice questions (understandability) and performing 3 perfective maintenance tasks (modifiability).</p>	<p>Apparently, diagrams with a high level of detail helped participants to understand the system better while those with a low level of detail helped to maintain it. Nevertheless, these results were fairly inconsistent and had no clear tendency.</p>

software defects, an improvement to the quality and a reduction in the software maintenance effort.

Garousi et al. (2013) present a multi-method empirical approach, and include a survey, a case study and some action-research in their proposal. A company providing satellite navigation system products was changing its software development processes and intended to measure which key factors were impacting on documentation usage (information sources, life-cycle phase, document type, roles, degree of experience, and patterns of usage) and which attributes were affecting the quality of the documentation. An exploratory survey was conducted with 25 employees to assess the aforementioned factors and attributes. Its results were reused in a case study during which the company data was employed, and the conclusion was reached that the usage of documentation differed according to the purpose for which it was used, e.g., documentation was more frequently used for development purposes than for maintenance purposes. Moreover, documentation the up-to-date-ness, accuracy and completeness of document artifacts were identified as the most important and relevant quality attributes as regards improving documentation efficiency. Finally, all these results were used in a set of action-research cycles for a continuous improvement of the company's documentation efficiency.

Finally, Fernández-Sáez et al. (2015) present the findings of a survey on the use of UML in software maintenance. A total of 178 professionals from 12 different countries took part in this survey. The main objectives of this survey can be summarized as follows: (i) to explore the extent to which UML diagrams are actually used in industry (59% of the answers indicated the use of a graphical notation, 43% UML), (ii) to acknowledge which was the most effective UML diagram for software maintenance (as expected, class, use case, sequence and activity diagrams), (iii) to discover what the perceived benefits of using UML were (less time needed for a better understanding and, thus, an improved defect detection), and (iv) to contextualize what kind of companies used UML documentation during software maintenance (larger teams seem to use UML more frequently).

To conclude, only a few evaluations of the benefits derived from the use of UML throughout the whole software development life cycle have been reported (Anda et al. 2006). This lack is even more evident in the software maintenance phase with regard to the benefits of employing UML models in source-code comprehensibility and modifiability. In addition, results sometimes seem to be contradictory or do not confirm the general belief, i.e., the UML documentation is not always beneficial to source-code comprehension/maintenance. In a few cases, it has been observed that UML documentation does not provide an advantage as regards time, while it helps to improve the correctness and quality achieved when solving complex tasks (e.g., Arisholm et al. 2006). In some other cases, UML is beneficial to source-code maintenance only if developers have a higher ability and/or the models have some specific characteristics (e.g., Fernández-Sáez 2015, 2016, Leotta et al. 2013). To try making things clearer, we present the results of long-term research conducted to study whether or not UML-based models (produced in both the requirements and analysis phases) aid source-code comprehensibility. In particular, we aggregate and synthesize the outcomes of 12 controlled experiments in different experimental contexts and on different sites with participants who had different levels of expertise. A total of 333 observations were obtained from these experiments. At the best of our knowledge, our study is currently the largest in the literature concerning the UML and its effect on source-code maintenance.

2.2 Model-based traceability

Traceability can be considered as an important related concept, since the comprehension of source code supported by the use of models is affected by the possibility of implic-

itly identifying relationships between source code and software models. When traceability information is explicitly documented in addition to the models, it can help developers to comprehend source code. In the following, we present related work focussed on this matter. For example, Lehnert et al. (2013) combine impact analysis, multi-perspective modeling, and horizontal traceability analysis to support the specification of models and the development of source code and test cases. They propose a unified meta-model approach that can be supplied by the Eclipse Modeling Framework (EMF) (2012) and a centralized model repository. The approach makes it possible to analyze the dependencies between software artifacts according to the type of change which is applied to them. The idea is to verify the interplay of change operations and dependency relations between models and code with the aim of identifying the propagation of further changes.

Hammad et al. (2011) propose an approach with which to automatically determine whether a change in the source code affects the design of the system (i.e., UML class diagrams). The aim of this approach is to maintain consistency between developed code and models by exploiting code-to-design traceability when the source code evolves. The proposed approach, along with the prototype implemented, was assessed by performing a case study based on the commits extracted from four open source projects during a three years period. The results revealed that most of the code changes do not impact on the software models. Furthermore, these commits regard a smaller number of changed files and less lines of code with regard to commits that impact on software models. Another interesting result is that most bug fixes do not impact design.

Settimi et al. (2004) assess the effectiveness of information retrieval techniques as regards tracing new and changed requirements to UML artifacts, code, and test cases. The authors summarize the most important result from their research as follows: tracing to UML elements provides a higher perspective of the proposed change than would be possible if links were generated directly to the code and supports the growing trend toward Model Driven Development. One possible implication is that this kind of link might reduce the effort required to analyze the impact of the changes.

Cariou et al. (2002) present an approach that focuses on object collaboration. This is recognized as an important building block for structuring object-oriented design in a distributed context. In an attempt to deal with the problem of the deterioration of the collaboration information during the detailed design process, a process and an architecture is proposed to preserve object collaboration information, from the analysis to design and implementation. The idea is to employ UML collaboration diagrams, with the addition of OCL constraints that follow specific rules to suit component specification requirements. This specification can be successively transformed into various low-level implementation designs depending on non-functional constraints by means of a refinement process.

Pavalkis et al. (2013) extend UML by defining a model-driven domain-specific language engine in order to manage traceability schemas and traceability analysis means. The authors specifically propose a framework that can be used to derive properties in order to trace project artifacts. The authors run several case studies to show how the framework can be used to adapt their solution to a particular development method and domain-specific language in a development process, and to automate the maintenance of traceability relations.

The proposal by Tang et al. (2007) is focused on the understanding of design rationale to support the detection of inconsistencies, omissions, and conflicts in an architecture design. The model incorporates design rationale, design objects and their relationships, and traceability methods are applied so as to change impact analysis and root cause analysis. The UML notation is used to represent the AREL model, which is an acyclic graph that relates

elements of the architecture to their rationale by exploiting the ARtrace directional link (namely, UML stereotyped association) (Tang et al. 2007).

Pavalkis et al. (2013) propose an approach for improving vertical traceability of UML models, thus eliminating the additional complexity involved in defining and maintaining traceability information in the software projects. Indeed, traceability information is not statically managed and memorized, but the use of derived properties allows the dynamic calculation of this information, which is then analyzed using dedicated and already existing tool-specific means. The application of the approach to a particular development process has shown that it allows the completeness of the project to be validated and the impact of changes to be analyzed, without affecting issues related to the management of traceability information.

A number of approaches has been proposed in the context of model-based traceability. Often the validity of these approaches has been empirically assessed through case studies, while the use of controlled experiments with participants seems marginal, especially to study the effect of model-based traceability on source-code maintenance.

3 Background

As the number of empirical studies grows, the need to aggregate evidence from multiple primary empirical studies (e.g., experiments) increases (Wohlin et al. 2012). There are two main reasons for aggregating evidence. Firstly, new research should always take existing knowledge into consideration as its starting point. That is, reviews summarizing the outcomes of various intervention trials are an efficient method by which to obtain the “bottom line” regarding what works and what does not. Secondly, primary empirical studies may together provide answers to research questions, when these studies alone are not sufficient to answer these research questions. In Section 3.1, we first briefly introduce strategies that can be used to summarize and synthesize outcomes from different empirical studies/experiments. When primary studies are synthesized using statistical methods (i.e., using a meta-analysis), it is crucial to verify whether or not these studies are heterogeneous (Pickard et al. 1998). In Sections 3.2 and 3.3, we present some background on how to deal with heterogeneity in meta-analysis studies. Finally, we present the process we have defined to deal with heterogeneity.

3.1 Aggregating results from primary studies

The collection, synthesis, and review of empirical evidence must comply with scientific standards. There are several strategies with which to summarize and synthesize outcomes from different empirical studies/experiments. For example, a systematic literature review is a means used to collect and synthesize empirical evidence from different empirical studies (Kitchenham and Charters 2007). A systematic literature review is referred to as a secondary study, while the empirical studies in such a review are referred to as primary studies. A systematic literature review has a research question, similar to that in primary studies. If the research question is more general (or if the field of research is less explored) a mapping study may be carried out.

When a set of empirical studies on a topic is collected, synthesis or aggregation takes place. A synthesis based on statistical methods is referred to as a meta-analysis (Wohlin et al. 2012). It can be applied to analyze the outcomes of several dependent and/or independent studies/experiments. The most important advantage of using a meta-analysis is that this kind

of secondary study makes it possible to achieve a higher statistical power for the variable of interest than primary studies. Although there is no accepted minimum number of primary studies in a meta-analysis, a minimum of 10 primary studies can be considered acceptable (Pickard et al. 1998).

3.2 Assessing heterogeneity

Studies may vary. In fact, the assumption that the studies are all representative samples of the overall true effect and only differ owing to sampling error is not always valid. In this situation, the studies are said to be heterogeneous (Pickard et al. 1998). It is useful to distinguish between different types of heterogeneity. According to Higgins and Green (2008), we can distinguish the following kinds of heterogeneity: clinical, methodological, and statistical. Variability in the participants, interventions, and outcomes studied may be described as clinical heterogeneity. Variability in study design and risk of bias may be described as methodological heterogeneity. Finally, variability in the intervention effects being evaluated in the different studies is known as statistical heterogeneity. This kind of heterogeneity can be considered a consequence of clinical heterogeneity or methodological heterogeneity, or both, among the primary studies. In the following part of this section we will focus on statistical heterogeneity and we refer to it simply as heterogeneity. This choice might affect the results of our investigation (see the discussion on the threats to validity in Section 5.4.4).

In a meta-analysis, the means usually employed to assess whether a set of single studies is heterogeneous the Cochran's Q test (Pickard et al. 1998). This test measures the deviation of observed effect sizes from an underlying overall effect size. The most frequently used cut-off point is 0.1. If the value is lower than this threshold, we can reject the null hypothesis (i.e., the primary studies are not heterogeneous) and we can then assume that studies are heterogeneous. The Cochran's Q test informs us only about the presence of heterogeneity, but it does not report on the extent of that heterogeneity (Huedo-Medina et al. 2006). Possible measures of heterogeneity are I-squared and tau-squared. The I-squared measure is the percentage of total variation across experiments that is owing to heterogeneity rather than chance. Thresholds for the interpretation of I-squared values can be misleading because the importance of inconsistency depends on two main factors: (i) magnitude and direction of effects and (ii) strength of evidence for heterogeneity (e.g., p-value from the chi-squared test, or a confidence interval for the I-squared measure). A guide to the interpretation of I-squared values is based on the intervals suggested by Higgins and Green (2008):

- 0% to 40%: heterogeneity might not be important;
- 30% to 60%: a moderate heterogeneity may be present among the primary studies;
- 50% to 90%: a substantial heterogeneity may be present among the primary studies;
- 75% to 100%: considerable heterogeneity may be present among the primary studies.

Tau-squared is an absolute measure of heterogeneity. It is a measure of the standard deviation of effect sizes across the experiments. Values greater than 1 indicate that primary studies are heterogeneous (Pickard et al. 1998).

Huedo-Medina et al. (2006) stated that I squared should be used as a complement to the Cochran's Q test. Since both Tau-squared and I-squared are measures of heterogeneity, these measures can be considered both as a complement to the Cochran's Q test. That is, the heterogeneity of the primary studies can be measured by either (or both) Tau-squared or I-squared if the Cochran's Q test rejects the null hypothesis that these studies are not heterogeneous.

If studies are not heterogeneous (i.e., they are homogeneous), they should be combined in a meta-analysis using a fixed effect model. This model assumes that the size of the treatment effect is the same (fixed) across all the experiments.

3.3 Dealing with heterogeneity

A number of options are available if (statistical) heterogeneity is identified among a group of primary studies that would otherwise be considered suitable for a meta-analysis. For example, Higgins and Green (2008) suggested:

1. *Check again that the data are correct.* Severe heterogeneity can indicate that data have been incorrectly extracted and/or used. For example, if standard errors have mistakenly been entered as standard deviations for continuous outcomes (Higgins and Green 2008).
2. *Do not do a meta-analysis.* If there is considerable variation in results, and particularly if there is inconsistency in the direction of the effect, it may be misleading to quote an average value for the effect.
3. *Explore heterogeneity.* The goal is to determine the causes of heterogeneity. This is problematic since there are often many characteristics that vary across primary studies from which one may choose. Heterogeneity may be explored by conducting sub-group analyzes. Each of these groups should have a minimum of 4 studies/experiments (Fu et al. 2011). Explorations of heterogeneity can at best lead to the generation of hypotheses. They should be interpreted with caution. Investigations of heterogeneity when there are very few studies are of questionable value.
4. *Ignore heterogeneity.* Heterogeneity can be ignored and a fixed effect meta-analyzes can be performed. The pooled effect estimate from a fixed effect meta-analysis is normally interpreted as being the best estimate of the intervention effect. However, the existence of heterogeneity suggests that there may not be a single intervention effect but rather a distribution of intervention effects. The pooled fixed effect estimate may be thus an intervention effect that does not actually exist in any population, and, therefore, have a confidence interval that is both meaningless and too narrow.
5. *Perform a random effects meta-analysis.* A random effects meta-analysis may be used to incorporate heterogeneity among primary studies. This is not a substitute for a thorough investigation of heterogeneity. It should be intended primarily for heterogeneity that cannot be explained.
6. *Change the effect measure.* Heterogeneity may be an artificial consequence of an inappropriate choice of the effect measure (dependent variable). Furthermore, the choice of the effect measure for dichotomous outcomes (odds ratio, relative risk, or risk difference) may affect the degree of heterogeneity among results. When control group risks vary, homogeneous odds ratios (or risk ratios) could lead to heterogeneous risk differences (and vice versa).
7. *Exclude studies.* Heterogeneity may be owing to the presence of one or two outlying studies. That is, the results of these studies could conflict with the results of the remaining studies. In general, it is unwise to exclude studies on the basis of their results. This may introduce bias into the meta-analysis results. However, if an obvious reason for the outlying result is apparent, the study might be removed with more confidence. It is advisable to perform analyzes both with and without outlying studies as part of a sensitivity analysis. Whenever possible, potential sources of diversity that might lead to heterogeneity should be specified in the experimental protocol. When excluding studies, a researcher should also take into consideration the number of studies to be then considered in the meta-analysis, thus avoiding losing the representativeness of results.

3.4 A process based approach to deal with heterogeneity

We employed the strategies by Higgins and Green (2008) as a basis to define a process with which to better deal with (devised statistical) heterogeneity in meta-analysis studies. Figure 1 shows this process as an activity diagram with object flow, where the activities are the phases of the process and the objects represent the input/output to these phases. In particular, the process first suggests checking data from individual experiments to be sure that they are correct (*CheckingData*). This is performed before checking that the secondary studies are heterogeneous. The phase in charge of assessing whether secondary studies are statistically heterogeneous is *AssessingHeterogeneity*. If the secondary studies are not heterogeneous, a meta-analysis will take place in *PerformingFixedEffectModel*. If the secondary studies are heterogeneous, the researcher can decide not to carry out a meta-analysis (the flow is placed in the end-node) or to go ahead with the process. In the latter case, the researcher can decide whether or not to ignore heterogeneity. Ignoring heterogeneity implies the execution of *PerformingFixedEffectModel*. On the other hand, the researcher can decide whether or not to incorporate the heterogeneity into a meta-analysis. If the researcher decides to incorporate the heterogeneity, the meta-analysis will take place by executing *PerformingRandomEffectsModel*. That is, a random effects model will be applied to all the primary studies. It is also possible to either explore or not explore heterogeneity. If heterogeneity is not explored, the researcher can decide to choose a different effect measure (or dependent variable) in order to aggregate the results in a meta-analysis (*ChangingEffectMeasure*). There are two ways in which to carry out an exploring: by excluding one or two outlying primary studies (*ExcludingStudies*) or by identifying sub-groups of experiments (*IdentifyingSubGroups*). In the case of the researcher excluding one or two outlying studies, meta-analysis takes place with a single group of studies. This implies that the process in Fig. 1 is instantiated only once. If sub-groups of experiments are identified, the process is instantiated for each of these sub-groups.

In order to show the different instances¹ of the process shown in Fig. 1, we used regular expressions in which the symbols are those used to label the phases of the process shown in Fig. 1: *a*, *b*, *c*, *d*, *e*, *f*, and *g*. For example, the label for *CheckingData* is *a*. We also considered the empty symbol (i.e., ϵ). The regular expression defined is shown as follows:

$$ab((eab)^*(gb)^*(fb)^*)^*(c|d|\epsilon)$$

We used regular expressions to provide a more compact representation of the instances of this process. In particular, our solution makes it possible to establish a clear link between the instances in our process and the sentences in the language described by means of the regular expression shown previously. We shall also use the sentences from the regular expression to facilitate our discussion on how we dealt with heterogeneity in the study presented in this paper. For example, the sentence *abc* means that we checked data, assessed heterogeneity, and performed a fixed effect model to aggregate results. However, our compact representation does not make it evident whether we execute *PerformingFixed EffectModel* because heterogeneity is ignored or because the experiments are homogeneous. We deal with this ambiguity by underlining sentences in the regular expression defined to indicate that the experiments are homogeneous. The sentence *abc*, therefore, indicates that the experiments are heterogeneous, while *abc* indicates that they are homogeneous. In both cases the

¹An instance of a process is a sequence of the activities/phases.

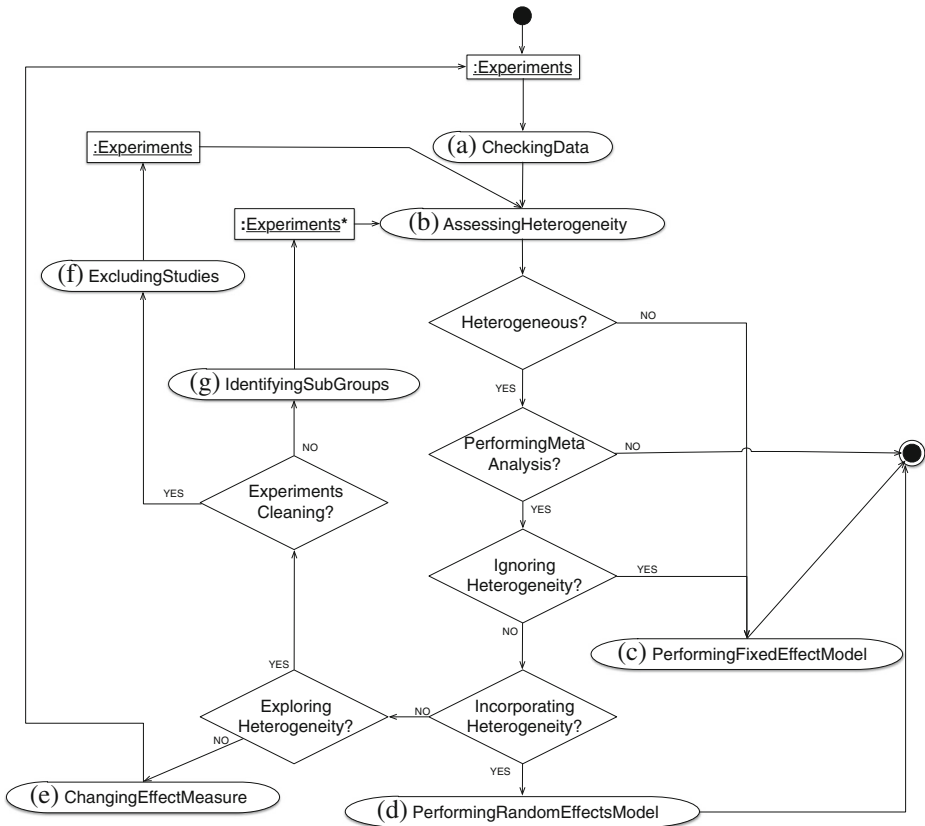


Fig. 1 The process defined to deal with statistical heterogeneity, in which each phase has a label (i.e., a, b, c, d, e, f, or g) associated with it. This allows a rapid reference to the phases in the process

sequences in which the phases are performed is: *CheckingData*, *AssessingHeterogeneity*, and *PerformingFixed EffectModel*.

4 Our long-term investigation

In a survey we conducted in 2009 (Scanniello et al. 2010), the main results suggested that, in order to deal with software maintenance and evolution tasks, many of the companies interviewed use UML diagrams produced in the requirements analysis² phase and, in a few cases, those produced in the design³ phase. The most frequently used UML diagrams were: use case, class, and sequence diagrams. Another result of this survey was: maintenance operations were performed by practitioners with a few years of experience in software devel-

²It is also called requirements engineering process and it is the process of determining user expectations (i.e., requirements) as regards a new or modified product.

³It maps the requirements onto the software architecture that defines the components, their interfaces and behaviors. The design document describes a plan with which to implement the requirements.

opment and maintenance and maintenance operations were supported by models produced in the analysis and design phases. To this end, companies generally employ developers with a Bachelor's or Master's degree in Computer Science and with between 1 and 5 years' experience.

The main results of the aforementioned industrial survey were used as a basis to begin the long-term research shown in this paper. We quantitatively studied to what extent developers understand source code when they were provided with source code alone or with source code and UML software models together. Our research consisted of the following two main directions, which were carried out in parallel:

- In the first direction of our long-term research, we studied the comprehension of source code when it was complemented with analysis models (i.e., models produced in the requirements engineering process) based on the UML notation: use case diagrams, class diagrams, and sequence diagrams. In particular, use case diagrams and use cases (textual description of the use cases in the use case diagrams) were employed to represent functional requirements. These requirements are represented with successful use cases. The participants were also given exceptional and/or boundary conditions. Class diagrams were used to abstract the objects from the problem domain (i.e., the object or conceptual model), while sequence diagrams (also produced in the requirements engineering phase) were used to model the dynamic and/or functional behavior of the software (Bruegge and Dutoit 2003). First, we conducted a pilot study with Computer Science Bachelor's degree students at the University of Basilicata. The results of this pilot study were presented in Gravino et al. (2010) and can be summarized as follows: the use of analysis models does not significantly improve the comprehension of source-code. We successively conducted a family of four controlled experiments on this subject (Scanniello et al. 2014), the goal of which was to strengthen the findings obtained in the pilot study. The experiments were carried out with students and practitioners from Italy and Spain who had different abilities and levels of experience with the UML. The results attained indicated that UML analysis models did not appear to improve source-code comprehensibility, thus confirming the results from the pilot study.
- In the second direction, we conducted two kinds of controlled experiments. The main goal of the second direction of our research was to study the comprehension of source code when it was complemented with design models (i.e., models produced in the design phase) based on the UML. In particular, the two kinds of controlled experiments were conducted in parallel and pursued the following main goals:
 - (i) Assessing the potential benefits derived from the use of UML class and sequence diagrams (both produced in the design phase) as regards the comprehension of object-oriented source-code (Gravino et al. 2015). Two experiments with Computer Science Bachelor's and Master's degree students were, therefore, conducted. The data analysis revealed that those participants with more experience of the UML and computer programming (i.e., Master's degree students) benefitted from the use of UML models produced in the design phase (from here on UML design models, also).
 - (ii) Investigating whether providing source code with UML class diagrams used to graphically document design-pattern instances⁴ improves source-code

⁴Design-pattern instances can be seen as a micro-architecture that developers copy and adapt to their particular designs in order to solve the recurrent problem described by the design pattern (Bruegge and Dutoit 2003; Gamma et al. 1995).

comprehensibility. This implies that multiple diagrams refer to the same piece of software and each of them can be seen as an excerpt of the entire class diagram of that piece of software. This is the main difference between the aspects (i) and (ii). We conducted an experiment with Master's degree students (Scanniello et al. 2010). The control group of this experiment comprised students who were given source code alone without any reference to the design-pattern instances contained in it. We carried out four successive controlled experiments with participants who had different experience as regards programming and software modeling (i.e., Bachelor's, Master's, and PhD students and practitioners) (Scanniello et al. 2015). The effect of textually documented design-pattern instances was also studied. Data regarding this kind of documentation was clearly not considered in the study presented in this paper, since the UML was not used.

All 12 experiments (primary studies) briefly described above and summarized in Table 2 were carried out by following the recommendations provided in Juristo and Moreno (2001), Kitchenham et al. (2002), Wohlin et al. (2012). In this section, we summarize the planning and the operation phases of these experiments and provide the most salient information concerning the study presented in this paper.

The experiments are reported according to the guidelines suggested by Jedlitschka et al. (2008). For replication purposes, we have made the raw data regarding of all our experiments available on the web.⁵

4.1 Goal

According to the Goal Question Metrics template (Basili and Rombach 1988), the goal of the study presented in this paper is: *to analyze* the use of UML analysis and design models *for the purpose of* understanding their utility *with respect to* the comprehensibility of object-oriented source code *from the point of view of* the software engineer *in the context of* students and practitioners.

4.2 Context selection

We used different software systems and UML diagrams in our study. The systems used were those described in the fifth column of Table 2, while the diagrams are those shown in the second column of this table. All the experimental objects were desktop applications implemented in Java. We used the Music Shop⁶ and Theater Ticket Reservation⁷ applications. Moreover, their models were created in a course on Advanced Object-Oriented Programming (AOOP). The lecturer of this course was involved in neither the study presented here nor those shown in Table 2, which allowed us to mitigate possible threats to construct validity (Experimenters' Expectancies). We used the source code that the lecturer selected from

⁵www2.unibas.it/gscanniello/SourceCodeComprMetaAnalysis/data.xlsx

⁶This is a software system that is used to sell and manage CDs/DVDs in a music shop. The feature *search for a singer* was used in the experiments: the user inserts a string (e.g., the surname of the singer), and the system then searches for all the singers that satisfy the search criterion and shows them in a list of the associated information.

⁷This is a software system with which to book and buy theater tickets. The feature *buy a theater ticket* was used in the experiments: the system shows the list of the available tickets for a given theater and performance, and the user then chooses the ticket and inserts data about the spectator.

among the software systems developed by the students on the AOOOP course. We did not have any control over the selection process. However, we reviewed the documentation and models to find possible issues. It was not necessary make any considerable modifications. We only removed possible typographical errors and indented source code when appropriate. Source-code comments were removed to avoid their presence having any effect on the results. That is, the effect of source-code comments could have been confused (or could have interacted) with the main factor being studied (see Section 5.4.2). The students that developed the experimental objects did not participate in the experiments.

In order to deal with the threat to external validity, we also used open-source software. We selected a chunk (i.e., a vertical slice) of JHotDraw v5.1. This chunk included: 10 instances of design patterns in total — two instances for the State design pattern and one instance for the following design patterns: Adapter, Strategy, Decorator, Composite, Observer, Command, Template Method, and Prototype — and the design patterns considered were well-known and widely adopted (Gamma et al. 1995). We documented the design-pattern instances present in the source code using both the JHotDraw documentation and the PMART dataset (Guéhéneuc 2007). This allowed us to document both intentional and unintentional design-pattern instances. Although traceability links are important as regards identifying relationships between source code and software models (see Section 2.2), we avoided providing them to the participants. The rationale for this was that making this information available to a developer could affect source-code comprehensibility in an undesirable way, i.e., concealing the effect of the use of software models on source-code comprehension. However, this design choice poses an additional threat to external validity, since traceability links could be available (e.g., post-facto) to support program comprehension tasks in the software industry (Antoniol et al. 2002).

We conducted all the experiments, with the exception of DePra (see Table 2), in research laboratories. DePra was conducted at the participants' companies. All the experiments were conducted under controlled conditions. The most participants' salient characteristics are summarized in Table 2 (third column). Participation in the experiments was on a voluntary basis. The participants were not paid. Each participant took part in only one experiment. Further information on the experimental objects and their selection process, along with the characteristics of the participants in the experiments, can be found in Gravino et al. (2015), Gravino et al. (2010), Scanniello et al. (2014), Scanniello et al. (2010), Scanniello et al. (2015).

4.3 Selection of variables

In each experiment, we considered those participants who were given source code alone as comprising the *control group*, while the *treatment group* comprised students who were given source code with software models based on the UML. The independent variable (from here on manipulated factor or main factor, also) considered in each primary study was, therefore, *method*. This variable is nominal and can assume the following two values: *models* (UML-based models and source code without comments) and *source code* (source code without comments).

The effect of the manipulated factor was analyzed on the following chosen constructs:

- **Comprehension.** This denotes the comprehension level of the source code achieved by a software engineer.
- **Completion time.** This denotes the time a software engineer takes to accomplish a comprehension task.

We used questionnaires to assess source-code comprehensibility. The correctness of the answers provided to these questionnaires was quantitatively evaluated by the F-measure, i.e., it was used to measure the comprehension variable and then estimate the comprehension construct. F-measure is computed as the overall balanced harmonic mean of precision and recall of the answers given to the questions. Answers were provided in the form of string items (e.g., a sequence of method/class names and/or the text messages shown to a user) and compared with the expected ones. F-measure values range in between 0 and 1. The higher the value of this variable, the greater the comprehensibility of source-code was.

We estimated the completion time construct using the overall time (expressed in minutes) taken to answer a comprehension questionnaire. The higher the value of time, the greater the effort⁸ required to accomplish a comprehension task.

4.4 Design

We used different kinds of designs in the experiments. As shown in Table 2, we used crossover designs in DeMscExp1, AnMscExp1, AnMscExp2, AnMscExp3, and AnPra. In the remaining experiments, we adopted the one-factor-with two treatments design (Wohlin et al. 2012). This kind of experimental does not suffer from the presence of a possible carry-over effect,⁹ while the crossover design may do so. It is worth noting that the design of some experiments was randomized, while in others we used the participants' ability as a blocking factor (see Table 2 for details). When applicable, randomization allowed us to mitigate carry-over that we had already analyzed in the primary studies. In all the experiments, the participants accomplished the task alone, that is, they did not work in a group to accomplish a comprehension task.

4.5 Experimental tasks and operation

All the participants were asked to fill in a comprehension questionnaire and a post-experiment survey questionnaire. The composition of both these questionnaires depended on the experiment and the tasks. We formulated the questions in the comprehension questionnaires using a similar form/schema. In addition, these questions were formulated to assess comprehension of the source code that we believed to be more relevant and concerned understanding concepts in this source code, which (as suggested by Sillitto et al. 2008) involved multiple relationships and software entities. Further details on the experimental tasks and the experimental procedure can be found in Gravino et al. (2015), Gravino et al. (2010), Scanniello et al. (2014), Scanniello et al. (2010), Scanniello et al. (2015).

4.6 Analysis procedure

Th results of a meta-analysis are commonly displayed graphically as “forest plots” (Ried 2008). This kind of pictorial representation provides a quick and easy means to illustrate the relative strength of treatment effects. Forest plots display point estimates and confidence

⁸The time was an approximation of comprehension effort. This complies with the ISO/IEC 9126 standard ISO (1991) (and subsequent versions), in which effort is the productive time associated with a specific project task.

⁹If a participant is tested first under the experimental condition A and then under the experimental condition B, she/he could potentially perform better or worse under condition B.

Table 2 Summary of the experiments*

Experiment	UML diagrams	Number of participants and kind	Design	Exp. Objects	Results	
					Comprehension	Completion time
AnBsc (Gravino et al. 2010)	Use case, class, and sequence	16 3rd year Bachelor's degree students	- One-factor- with more treatments - Randomized	A chunk of Music Shop software system implemented in Java A chunk of a Theater Reservation system implemented in Java	The difference is not statistically significant	Not considered
DeBscExp1 (Gravino et al. 2015)	Class and sequence	16 2nd year Bachelor's degree students	- Crossover design- Randomized		The difference is not statistically significant	The difference is statistically significant. More time taken when models used
DeMscExp1 (Gravino et al. 2015)		16 second year Master's degree students			The difference is statistically significant. Better comprehension when models used	The difference is not statistically significant

Table 2 (continued)

AnMscExp1 (Scanniello et al. 2014)	Use case, class, and sequence	24 Master's degree students	- Crossover design - Ability as blocking factor	The difference is not statistically significant	Not considered
AnMscExp2 (Scanniello et al. 2014)		22 Master's degree students		The difference is not statistically significant	
AnMscExp3 (Scanniello et al. 2014)		22 Master's degree students		The difference is statistically significant. Better comprehension when models not used	
AnPra (Scanniello et al. 2014)		18 Practitioners		The difference is not statistically significant	
DeMscExp2 (Scanniello et al. 2010)	Class	24 Master's degree students	-One-factor with more treatments - Ability as blocking factor for students - Years of working experience as blocking factor for practitioners	Analysis of comprehension not presented	The difference is statistically significant. Less time taken when models used
			JHotDraw: a two-dimensional graphics framework for structured drawing editors implemented in Java. A chunk (i.e., vertical slice) of the entire software was selected		

Table 2 (continued)

DePra (Scanniello et al. 2015)	16 Practitioners	The difference is statistically significant. Better comprehension when models used	The difference is not statistically significant
DeMscExp3 (Scanniello et al. 2015)	16 1st year Master's degree students	The difference is not statistically significant	
DeBscExp2 (Scanniello et al. 2015)	15 3rd year Bachelor's degree students		
DePhd (Scanniello et al. 2015)	10 Ph.D students	The difference is statistically significant. Better comprehension when models used	

* The labels of the experiments are expressed using the upper-camel case notation. Each compound word has a specific meaning. The first word suggests the kind of UML-based documentation: (An) stands for those documents produced in the analysis phase and (De) stands for those documents produced in the design phase. The second word indicates the kind of participant. For example, Pra and Msc stand for practitioner and Master's degree student, respectively. The third word is a progressive id for the experiment. We used id to better discern each experiment that was based on the same kind of documentation and involved the same kind of participants.

Table 3 Instances of the process shown in Fig. 1 considered

Sentence	Dependent variable
<i>ab</i>	Comprehension, Completion Time
<i>abc</i>	Comprehension, Completion Time
<i>abd</i>	Comprehension, Completion Time
<i>abgbc</i>	Comprehension, Completion Time
<i>abgbd</i>	Completion Time
<i>abfb</i>	Comprehension, Completion Time
<i>abfbc</i>	Comprehension, Completion Time
<i>abfbd</i>	Comprehension, Completion Time
<i>abeab</i>	Efficiency
<i>abeabc</i>	Efficiency
<i>abeabd</i>	Efficiency

intervals for individual experiments, in addition to an estimate of the overall summary effect size. This notation also shows the extent to which each experiment contributes to the overall result.

5 Results and discussion

In Table 3, we show some sentences from the language defined by employing the regular expression reported in Section 3.4. For example, *ab* covers the case of no heterogeneity in which a meta-analysis is not executed, while *abc* and *abd* represent the case in which heterogeneity is ignored or incorporated and a meta-analysis is executed, respectively. Conversely, in order to cover the case of exploring heterogeneity, we performed sub-group analyzes for the dependent variables Comprehension and Completion time, since the experiments were heterogeneous. In particular, we grouped the experiments according to the kind of models used, namely An (i.e., UML-based models produced in the requirements engineering process) and De (UML-based models produced in the design phase). These cases are covered by the sentences $ab(gb)^*c$ and $ab(gb)^*d$.

Furthermore, $ab(fb)^*$, $ab(fb)^*c$, and $ab(fb)^*d$ cover the cases of the application of the process shown in Fig. 1 when we cleaned the experiment set by excluding those studies that involved participants with little experience, namely AnBsc, DeBscExp1, and DeBscExp2. We also considered the cases of not exploring the heterogeneity and employing a further dependent variable, namely Efficiency,¹⁰ which are covered by the sentences $ab(eab)^*$, $ab(eab)^*c$, and $ab(eab)^*d$.

In the following subsections, we first report the results obtained and then discuss them according to the cases shown in Table 3. We then present the implications of our study and conclude with a discussion regarding the threats to validity.

¹⁰Efficiency is a derived measure that is computed as the ratio between comprehension and completion time. Task efficiency is a ratio measure and estimates a participant's efficiency as regards the execution of a comprehension task. The larger the efficiency value, the better it is. The perspective we adopted is that of quality in use (e.g., ISO 2000, 2011), since efficiency measures source-code comprehension achieved during the expenditure of available models.

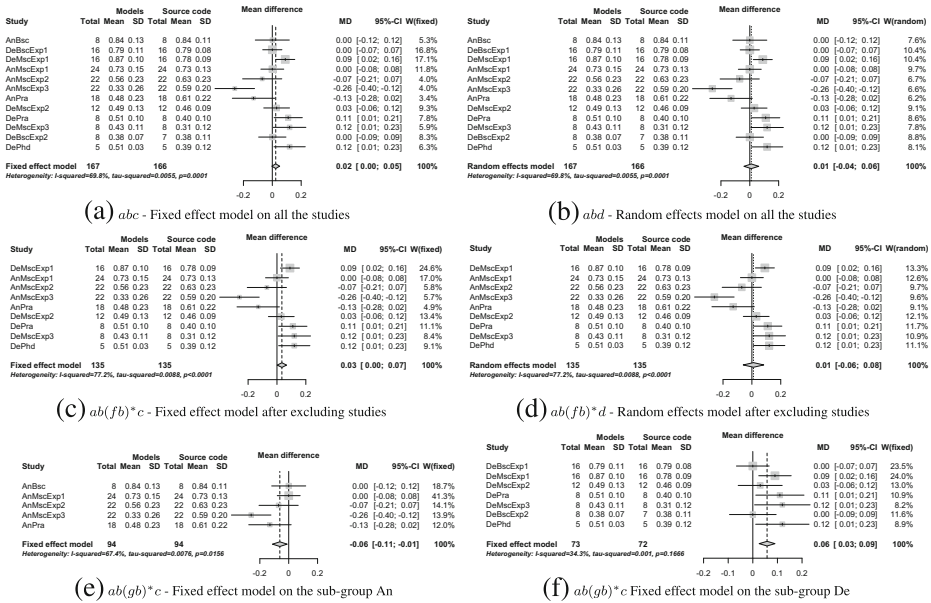


Fig. 2 Forest plots for Comprehension

5.1 Meta-analysis results

We summarize the results of each experiment, by employing the descriptive statistics of the measures of the dependent variables. The descriptive statistics for Comprehension (i.e., mean, standard deviation, and number of observations) grouped by method are shown in the forest plot in Fig. 2 (for each of the cases considered shown in Table 3). The same descriptive statistics for Completion time and Efficiency are reported in Figs. 3 and 4, respectively. It is worth mentioning that we do not show any results for the sentences *ab*, *abfb*, *abgb*, and *abeab* (see Table 3), because the process in Figure 1 does not require the execution of a meta-analysis. Some other sentences are also not reported (e.g., those that excessively reduce the number of experiments in the analyzes, see Section 3).

The results are synthesized by means of the Mean Differences (MDs) of the outcome measures of the experiments. This is possible because the experiments have the same outcome measures for each dependent variable studied.

Results concerning the testing of heterogeneity are also shown at the bottom of each forest plot (see, for example, the left-hand side of Fig. 2a). With regard to Comprehension, the results of the Cochran’s Q test (see Fig. 2a — sentence *abc*) suggest that the experiments were heterogeneous ($p=0.0001$) while the I-squared values indicates a substantial/considerable heterogeneity that is also confirmed by the Tau-squared value. We can, therefore, incorporate heterogeneity and apply a random effects model (see Fig. 2b — sentence *abd*).

As the squares in the figure suggest, the experiments contributed equally to the overall result, that is, the use of models slightly improved source-code comprehensibility (MD = 0.01). Indeed, source-code comprehensibility was not significantly different¹¹ when using

¹¹Effect size is statistically different from the overall effect if the diamond (at the bottom of the forest plot) does not intersect the vertical line.

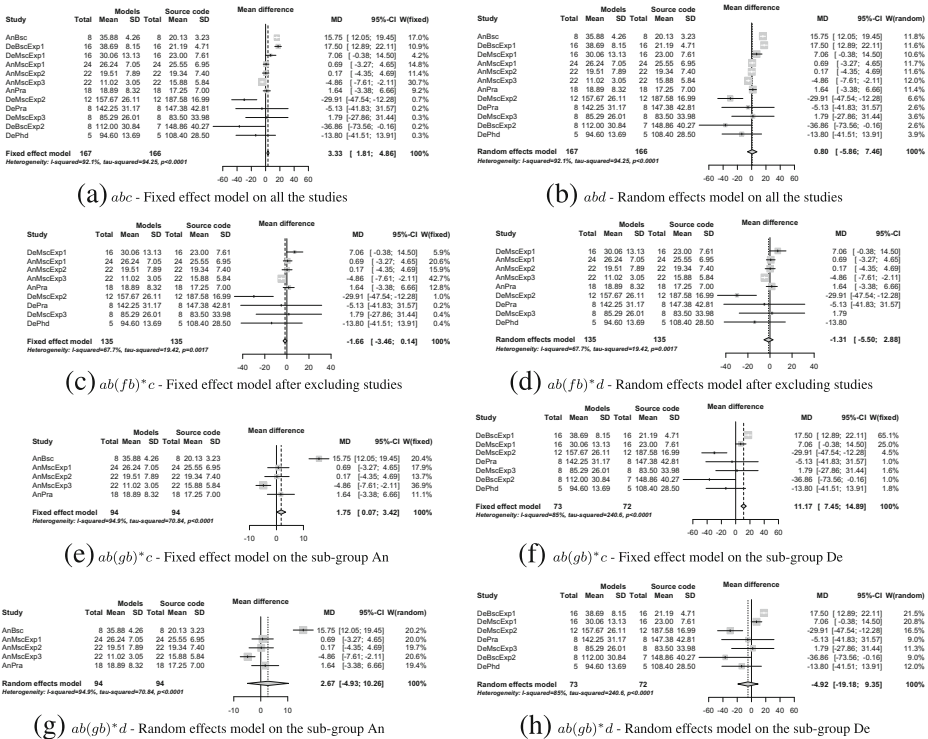
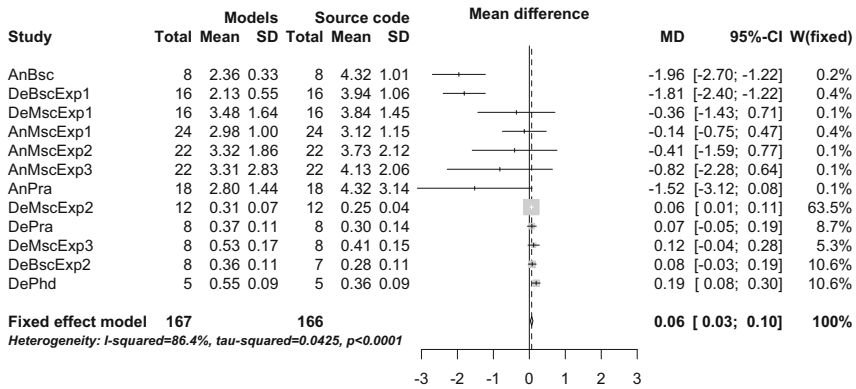


Fig. 3 Forest plots for Completion Time

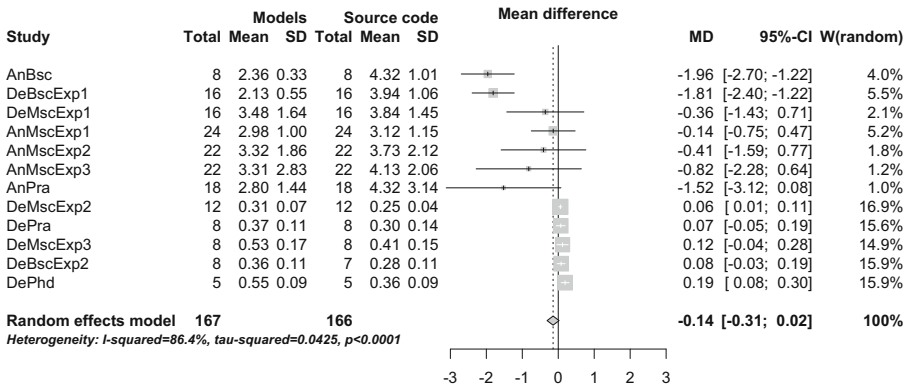
or not using models in the execution of comprehension tasks. This result is also confirmed by the overall 95% confidence interval¹² (IC) whose value is $[-0.04; 0.06]$. If we decide to ignore heterogeneity and apply a fixed effects model, we obtain the MD values shown in Fig. 2a. As we can see, the result is quite similar when taking into account the MD values. However, the IC value is $[0.00; 0.05]$.

With regard to Completion time, in Fig. 3a (sentence *abc*), the experiments, according to the Cochran’s Q test, are heterogeneous ($p < 0.0001$). By incorporating heterogeneity and applying a random effects model we obtain an MD value of 0.8 and an IC value of $[-5.86; 7.46]$. Unlike Comprehension, the squares are not proportional in size when studying Completion time (see Fig. 3b — sentence *abd*). This signifies that some experiments contributed to the overall result more than others. The forest plot suggests that the difference in the completion time is not significant when using or not using models to accomplish comprehension tasks. However, when ignoring heterogeneity and applying a fixed effects model (see Fig. 3a — sentence *abc*), the completion time is statistically different when using or not using models in comprehension tasks. The MD value is 3.33 while the IC value is $[1.81; 4.86]$. The choice of how to manage heterogeneity, therefore, proves that the results are crucial in this case, and influences the overall results regarding the impact of UML models on source-code comprehensibility.

¹²This is a range of values that we are 95% certain that it contains the true mean value.



(a) *abeabc* - Fixed effect model obtained when changing the effect measure and considering all the studies



(b) *abeabd* - Random effects model obtained when changing the effect measure and considering all the studies

Fig. 4 Forest plots for Efficiency

Since our experiments were heterogeneous with regard to the two dependent variables, we also decided to explore heterogeneity and perform sub-group analyzes. In our study, for both the dependent variables introduced in the design of the study (see Section 4.3), we can group experiments according to the kind of models: An and De. We postulated (on the basis of the results of the primary studies) that models produced in the design phase are closer to source code than those produced in the analysis phase. In other words, we could expect that De would aid source-code comprehensibility, while An would not. The forest plot for An and comprehension is shown in Fig. 2e. We used a fixed effects model, for the sentence *abgbc*, because of the results of the heterogeneously analysis. That is, it can be considered that the studies are not heterogeneous since the Cochran’s Q test returned 0.0156 as the value for *p* (this is why the the sentence of the regular expression aforementioned is underlined). The MD value obtained is low (-0.06) and the IC value is [-0.11; -0.01]. It would appear that the presence of UML analysis models does not aid source-code comprehensibility. The assumption made in order to explore heterogeneity and stated above was, therefore in some respects confirmed. The observed results thus allow us to state that UML analysis models

focus on the problem domain of the software (the environment in which the software will work) and do not provide any support as regards performing comprehension tasks on source code. Indeed, this kind of models could have confused the participants while comprehending the source code. For example, it could be possible that the participants trusted the models and did not pay adequate attention to the source code. These results are, perhaps, not overly surprising, but they are acceptable, as evidence/postulations need to be empirically verified and/or reaffirmed through the use of empirical studies (Basili et al. 1999; Kitchenham et al. 2002; Shull et al. 2008).

Figure 2f shows the forest plot of *abgbc* for the variable Comprehension in the experiments in which the participants were provided with design models (De). The experiments are not heterogeneous since the Cochran's Q test returned a value greater than 0.1 (i.e., $p = 0.1666$). We applied a fixed effects model, since the experiments can be considered homogeneous and it is for this reason that the sentence of the regular expression is underlined. Some experiments contributed to the overall result more than others. The most remarkable outcome is that the difference between using or not using models is significant. Those participants provided with design models understood the source code better because the diamond is on the right-hand side of the vertical line. The MD value is sufficiently large (i.e., 0.06). It is worth mentioning that only for the sentence *abgbc* for the variable Comprehension and the groups An and De, we observed an homogeneity of the experiments and applied a fixed effects model (i.e., *abgbc*).

With regard to Completion time, the meta-analysis results for the An and De sub-groups (sentences *abgbc* and *abgbd*) are summarized in Fig. 3e and f and g and h, respectively. Note that the experiments were heterogeneous (see the Cochran's Q test, which returned values less than 0.1 in all the cases). It was for this reason that we applied both the fixed- and random- effects models, i.e., we ignored and incorporated heterogeneity, respectively. The plots shown in Fig. 3g and h suggest that the time taken to complete a task was not significantly different when using or not using An and De models and exploiting a random-effects model. More specifically, and on the basis of the descriptive statistics, we can deduce that the participants provided with analysis models needed slightly more time to accomplish a comprehension task. Those in the other experiments spent less time when accomplishing the task with design models. In other words, it would appear that the use of design models paid back the time needed to read them because the effort required to comprehend source code decreased when compared with the effort of the participants provided with only source code. It is also worth mentioning that the exploration of the heterogeneity indicated that the groups of experiments in An and De are not heterogeneous for Comprehension, while they are heterogeneous for Completion time. This confirms that we have presented in Section 3.3, i.e., heterogeneity is not only linked to the kind and type of primary studies but also to the effect measure.

When exploring heterogeneity, one alternative to sub-group analysis is that of carrying out experiment cleaning. We considered the participants' experience in order to exclude experiments. We excluded AnBsc, DeBscExp1, and DeBscExp2 because the participants in these experiments were Bachelor's degree students. The forest plots for Comprehension shown in Fig. 2c and d show that the Cochran's Q test suggests that the experiments were heterogeneous since the value of p is less than 0.0001 (the Tau-squared and I-squared values indicated a good extent of such heterogeneity). That is, the experiments were still heterogeneous after this cleaning. It would appear that the participants' experience was not a cause of heterogeneity. We, therefore, incorporated heterogeneity and applied a random effects model (i.e., that shown in Fig. 2d). It is easy to observe that the model is quite similar to that shown in Fig. 2b, thus confirming the analysis performed when deciding to

incorporate heterogeneity. Indeed, source-code comprehensibility was not significantly different when using or not using models, and the MD value obtained is 0.01 while the IC value is $[-0.06; 0.08]$.

Similar to that which occurred with Comprehension, the forest plots for Completion time (see Fig. 3c and d) suggest that we can consider the experiments to be heterogeneous. Indeed, the Cochran's Q test suggests that the experiments were heterogeneous ($p < 0.1$) and the Tau-squared and I-squared values indicated a good extent of this heterogeneity. We, then, incorporated heterogeneity and applied a random effects model. Figure 3d shows that the squares are not proportional in size, i.e., some experiments contributed to the overall result more than others. Moreover, the forest plot suggests that the difference in the completion time is not significant when using or not using models, and the MD value obtained is -1.31 while the IC value is $[-5.50; 2.88]$.

The process shown in Fig. 1 was also used to analyze the case of: (i) not exploring heterogeneity and (ii) changing the effect measure by exploiting Efficiency (see Table 3). The results of the Cochran's Q test ($p=0.0001$) shown in Fig. 4a and b suggest that the experiments were heterogeneous. The results of the I-squared indicated a considerable heterogeneity (86.4%). This contrasts with the results of the Tau-squared measure, which suggests that the groups of experiments are not heterogeneous. We considered that experiments were heterogeneous owing to the results of the Cochran's Q test. We, therefore, incorporated heterogeneity and applied a random-effects model. The results obtained are summarized in Fig. 4b. Note that the squares are not proportional in size, i.e., some experiments contributed to the overall result more than others. Moreover, the forest plot suggests that efficiency is not statistically different when using or not using models, and the MD value obtained is -0.14 while the IC value is $[-0.31; 0.02]$. Conversely, upon ignoring heterogeneity and applying a fixed effects model there is a statistically significant difference when using or not using models (see Fig. 4a). The MD value obtained is 0.06 and the IC value is $[0.03; 0.10]$. As in the case of *abc* and *abd* in Fig. 3, the choice of how to manage heterogeneity proves to be crucial and can influence the analysis of the impact of models.

In Table 4, we report a summary of the results shown before. For each dependent variable, we show the effect of the method according to each case considered in our analysis (i.e., some of the paths of the process shown in Fig. 1, illustrating the sentences of the regular expression defined in Section 3.4) and the corresponding MD and CI values obtained from the meta-analysis. We present this table to distill the results from our meta-analysis and to make them easier to understand and to support the discussion of the results in the next subsection as well.

5.2 Discussion

The participants in our long-term investigation obtained slightly better scores for comprehension when using analysis and design models (see row *abc* in Table 4). The effect of the method is not significant. We can, therefore, postulate that models do not help a lot participants to comprehend source code, because these models do provide additional information on the subject application. In addition, the participants spent more time comprehending source code. This could be related to the effort needed to infer the information provided by the models that was definitively not useful as regards attaining an improved comprehension of source code.

We further investigated this aspect by performing sub-group analyzes, excluding cases, and changing the dependent variable. The results obtained from the meta-analyzes of the

Table 4 Summary of results. Values for MD and CI are reported only in the case of a significant effect of the method

Sentence	Dependent variable	Effect of the method	MD	95%-CI
<i>abc</i>	Comprehension	No	-	-
<i>abd</i>	Comprehension	No	-	-
<i>abgbc</i>	Comprehension	Yes for both An and De	-0.06; 0.06	[-0.11; -0.01]; [0.03; 0.09]
<i>abfbd</i>	Comprehension	No	-	-
<i>abc</i>	Completion time	Yes	3.33	[1.81; 4.86]
<i>abd</i>	Completion time	No	-	-
<i>abgbc</i>	Completion time	Yes for both An and De	1.75; 11.17	[0.07; 3.42]; [7.45; 14.89]
<i>abgbd</i>	Completion time	No	-	-
<i>abfbd</i>	Completion time	No	-	-
<i>abeabc</i>	Efficiency	Yes	0.06	[0.03; 0.1]
<i>abeabd</i>	Efficiency	No	-	-

two sub-groups suggest that the use of models affects source-code comprehensibility (see, for example, *abgbc* in Table 4), but in two slight different directions:

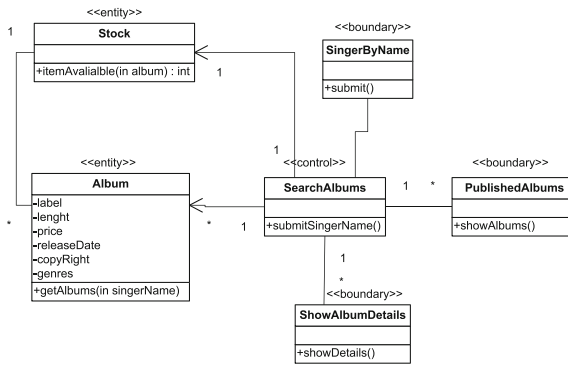
- the use of analysis models reduces source-code comprehensibility and increases the time taken to complete comprehension tasks;
- the use of design models improves source-code comprehensibility and reduces task completion time.

This outcome should not be at all surprising because analysis and design models are created for different purposes, although companies seem to ignore this difference (Scanniello et al. 2010) as we also discussed in the introductory part of Section 4. In particular, analysis models are created to capture a domain, to represent a set of requirements, to understand a poorly focused problem boundary, and so on, while design models can be used to structure source code and capture design artifacts that do not directly emerge from requirements (Bruegge and Dutoit 2003). As such, analysis models say little or nothing about source code and the use of this kind of models will not, therefore, benefit comprehension. As an example, the models of the Music Shop experimental object (i.e., that used in the experiments carried out by An group) shown in Fig. 5 do not provide implementation details. However, some details and some design decisions could be inferred. For example, some of the classes to be understood in Fig. 5b are in the source code (not reported here owing to their scant relevance, but available on the web for download) and also in the class diagram (available on the web for the download) in the experiments of the De group. That is, some classes in the problem domain are present in both the solution domain and the source code, thus possibly allowing the participant to also obtain a little information on the implementation from the analysis models. This scenario is customary for more traditional development approaches (e.g., Unified Software Development Process) (Bruegge and Dutoit 2003). With regard to an example of design decision, the architectural pattern employed could be used in an inferred manner. In particular, we could postulate that the architectural pattern implemented in Music Shop is the Model-View-Controller, given the division of the classes in: boundary, control, and entity.

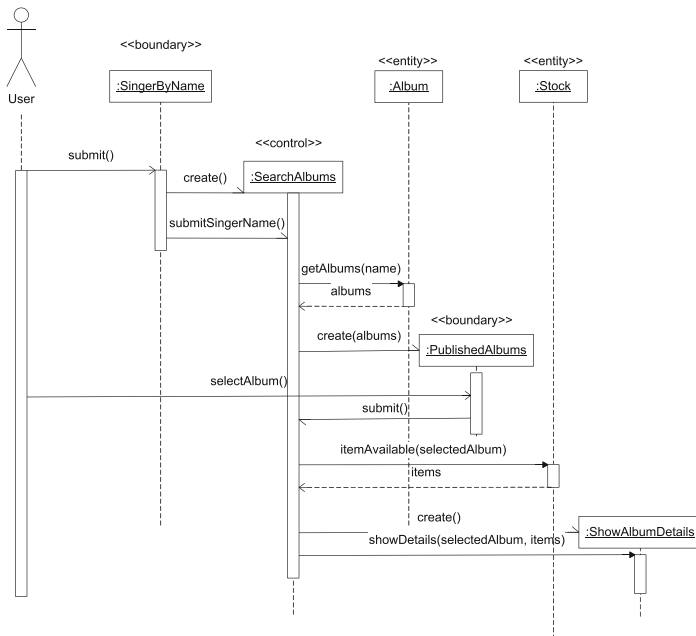
In both the cases mentioned above, it is the developers' ability (and possibly his/her knowledge of the subject software and its domain) that could make the difference in terms of source-code comprehension, rather than the actual information the analysis models provide.

Use Case Name: <i>Search Album by Singer</i>
ID: 3
Brief description of the use case: The user selects an album by specifying the singer's name. Details of the selected album will be shown by the system (e.g., original release date)
Main Actors: User
Flow of events: <ol style="list-style-type: none"> 1. The User inserts the singer name. 2. The system presents the album list of the specified singer. 3. The user selects the album of interest. 4. The system shows details on the chosen album and the number of available copies in the stock. As result, the presented information is: original release date, label, copyright, genres, length, price, and number of available copies.
Pre-condition: The user has selected the functionality search singer by name.
Post-condition: 1. The system shows details about a given album for the chosen singer.

(a) The *Search Album by Singer* Use Case



(b) Class diagram



(c) Sequence diagram for the functionality *Search Album by Singer*

Fig. 5 Some of the analysis models for the Music Shop application

In summary, we can assume that design models help more than analysis models in the comprehension of source code, since they are focused on the implementation aspects.

On the basis of our considerations, we can summarize our research results as follows:

- *Software models produced in the design phase aid in source-code comprehensibility.*

Despite the fact that our results improved the findings obtained for the individual experiments conducted in our long-term research, we cannot provide conclusive findings concerning whether analysis models helped in the understanding of source code in the context of graduate, undergraduate, PhD students, and novice practitioners.

5.3 Implications

We judged the implications of our study on the basis of the perspectives of practitioner/consultant (from here on simply practitioner) and researcher. When applicable, we also discuss future research directions related to these implications.

- UML-based modeling is important as it allows an improved comprehension of source code. These models should focus on aspects related to the solution domain (i.e., implementation aspects) of a subject software. Models produced in the analysis phase could, therefore, be considered of less importance if they are only intended to support the comprehension of source code. Furthermore, these models are of primary importance when they contain the subsequent development phases. We can speculate on this point because we used the same software in some of the experiments, but the models were produced in either the analysis or the design phases (e.g., AnBsc and DeMscExp1). From the practitioner's perspective, this result is relevant because it could be useful to adopt a development process based on the use of the UML. It might, however, be useless to give UML-based analysis models to the software engineer when he/she has to perform small maintenance operations on source code. This is aligned with the findings of Arisholm et al. (2006), Dzidek et al. (2008) (see Section 2) whose authors stated that the UML only seemed to be really useful as regards understanding complex systems. That is, this kind of models should be only used to support the subsequent phases of the development or to improve the comprehension of functional requirements (Abrahão et al. 2013). From the researcher's perspective, it would be interesting to investigate whether variations in the context (e.g., larger systems and more or less experienced software engineers) might lead to different results.
- UML analysis models appear to uselessly overload participants when performing comprehension tasks. Once again, the results obtained in our study coincide with those from some of the related work (e.g., Fernández-Sáez 2015, 2016) in which only design models appeared to help achieving a better understanding of the systems. This result is relevant for the researcher because it would be interesting to carry out further research into this aspect and discover in which context it holds.
- Although we are not sure whether our findings scale up to real projects, the results obtained could be true in all those cases in which the models are concerned with a part of the entire software and maintenance operations are performed on a chunk of the source code of the entire system.
- We observed that the models produced in the design phase aid the comprehension of source code and postulated that this is because they are closer to source code than those produced in the analysis phase. It might, therefore, be expected that reverse-engineered diagrams could also be at least as effective as the forward designed diagrams as regards aiding source-code comprehension since they are obtained directly from the source

code. Our results and those of Fernández-Sáez et al. presented in Fernández-Sáez et al. (2015) provide the basis for future work in this direction. This point is clearly relevant for any researchers who might be interested in studying the delineated research direction. If future work confirms that reverse-engineered diagrams are as effective as forward ones, the practitioner could be further motivated to use reverse engineering tools in his/her company.

- We focused on desktop applications. The models of these systems were sufficiently realistic for small-sized in-house software and subcontracting development projects. From the researcher's perspective, the effect of analysis and design models on different types of systems (e.g., smartphone and web apps) represents a possible future direction for our research. This point is clearly relevant for the researcher.
- The UML is widely used in the software industry. The results achieved are, therefore, useful for all the companies that exploit the UML as a support when software engineers are executing comprehension tasks (e.g., performing maintenance/evolution operations). Studies on this notation are currently required to understand the cases in which its use improves the comprehensibility and maintainability of source code. There are currently only a few evaluations, as we have discussed in the related work section.
- Dealing with heterogeneity could be crucial, since it may influence the results of the meta-analysis. This point is clearly of interest for any researchers who might be interested in studying how to deal with heterogeneity when integrating results from several studies by using a meta-analysis.

5.4 Threats to validity

Despite our efforts to mitigate as many threats to validity as possible, some are unavoidable. In order to comprehend the strengths and limitations of our empirical study, the threats that could have affected the results are presented and discussed as follows.

5.4.1 Internal Validity

This kind of validity is of concern when causal relations are examined. In our study, the design of the experiments could have affected the results. Different kinds of design were considered in each experiment. Each group of participants either worked on two different tasks with and without models (maturation or diffusion or imitation of treatments) or worked on a single task either with models or without models. The artifacts used to carry out the experiments (e.g., comprehension questionnaire and documentation) could also have negatively affected the experiments and thus the outcomes of the meta-analysis. We mitigated these threats by accurately designing all the material used in each experiment. In many cases, pilot studies were conducted to assess this experimental material. Threats to internal validity could also depend on how the participants are selected from a larger group. How the experiments were selected could also have affected the results. With regard to multiple group experiments the results could have been biased because of the different behavior of participants in different groups (i.e., interactions with selection). Social threats to internal validity could also have been present in the experiments.

5.4.2 External validity

Performing experiments with students could lead to doubts concerning their representativeness when compared to software professionals (interaction of selection and treatment). In

addition to experiments with students, we also carried out experiments with professionals and PhD students. It is worth mentioning that the tasks used did not require a high level of industrial experience. This led us to believe that the use of students was not a major issue here (Carver et al. 2003). However, if tasks are too simple they may be not representative. This could imply threats to the validity of the results.

Another threat to external validity concerns the experimental objects (interaction of setting and treatment). For example, we removed comments from source code and did not provide any explicit information on the traceability links between models and source code. We took these decisions to avoid the effect of source-code comments and traceability links being confused with the main factor studied.

With regard to source-code comments, some further considerations are required: (i) comments and source code may not be coherent (i.e., the comment does not describe the intent of the method and its actual implementation) with one another (Corazza et al. 2016) and (ii) it is possible that experienced developers (e.g., professionals) do not use comments (because they are often not updated when source code changes Fluri et al. 2007, Jiang and Hassan 2006) while performing comprehension tasks (Salviulo and Scanniello 2014). The first point could be dealt with by modifying the comments to make them coherent with the source code. Modifications to the comments could aim to restructure them to improve their readability (and then possibly comprehensibility). Any modification to source-code comments could affect external validity. As for the second point, we can do little or nothing. In fact, it could be that professional developers (and possibly PhD students) do not take much care with the comments, while those with little experience do. We, however, advise the use source-code comments and traceability links in future studies. Our research provides the basis for future work on this matter. We also indented the source code when preparing the experimental objects. This design choice could affect external validity. However, many of the available IDEs provide a feature to remove this kind of smell from source code. Therefore, it could happen that some of the participants indent the code, while other no. If this happened a series of threats to the conclusion validity could be risen. That is, if any effect of this kind of smell could be present it could affect results in an undesirable and uncontrollable way.

5.4.3 Construct validity

This kind of validity may have been influenced by the measures used to obtain a quantitative evaluation of source-code comprehensibility (inadequate preoperational explanation of constructs). Construct validity might also have been affected by the comprehension used and the post-experiment survey questionnaires, in addition to social threats. We employed post-experiment survey questionnaires designed using standard approaches and scales (Oppenheim 1992). The responses to this kind of questionnaire were used to explain the quantitative results. Another threat to external validity is mono-operation bias. All the experiments in our study included a single independent variable (or treatment). This may have under-represented the construct and thus not provided the full picture of the theory. The threat concerning the interaction of different treatments is not present in our research because the participants were involved in only one experiment. In order to mitigate construct validity, we conducted external replications and their results were subsequently aggregated with the other experiments and replications. The fact that the participants in our family of experiments received the documentation in their native language might also affect the validity of results.

5.4.4 Conclusion validity

This kind of validity concerns the ability to draw correct conclusions. In order to deal with conclusion validity, we performed a statistical analysis of the data gathered. Despite our effort, threats to low statistical power could still have been present. The number of experiments considered in our study might also have affected the results. With regard to the selection of the populations, we drew fair samples and conducted our experiments with participants belonging to these samples. Another threat to conclusion validity could be related to the number of participants. This kind of threat was mitigated because our study was based on 333 observations. The reliability of measures might affect results. In each individual experiment, the experimenters attempted to mitigate this kind of threat as much as possible. The threat related to the random heterogeneity of subjects could have been present in our meta-analysis study and in each experiment. We took this aspect into account in the meta-analysis. Finally, we did not deal with clinical heterogeneity or methodological heterogeneity. As mentioned in Section 3.2, statistical heterogeneity can be considered as a consequence of either, or both, clinical heterogeneity or methodological heterogeneity.

6 Conclusion

In this paper, we have presented the results of long-term research into the effect of using UML-based modeling in source-code comprehensibility. We removed comments from source code and did not provide any explicit information on the traceability links between models and source code. This research began in 2009 with an industrial survey (Scanniello et al. 2010), and the results of this survey were then used as a basis to conduct a number of controlled experiments (internal and external replications) with students and practitioners from Italy and Spain. The results of individual experiments were synthesized by means of a meta-analysis and presented in this paper. The most important outcome is: the use of UML models is important as regards allowing software engineers to better understand source code, given that these models focus on aspects related to objects (or entities) in the solution domain of a subject software and source code does not include comments. Models produced in the analysis phase are of less importance if their purpose is solely to enable the comprehension of source code.

Possible future directions for our research are: (i) performing further experimentation considering different and larger software systems related to unknown domains in order to verify whether the findings obtained are still valid; (ii) studying the effect of providing the participants with information in an incremental manner; (iii) analyzing the effect of different UML diagrams; (iv) investigating the effect of the same UML diagrams on non-source code comprehension tasks; and (v) studying the effect of the UML diagrams combined with source-code comments and traceability information on source-code comprehensibility in specific contexts (e.g., safety-critical software development).

Acknowledgements The authors would like to thank the participants in the experiments and all the people who supported the research presented in this paper. This work has been partially supported by the SEQUOIA Project, (TIN2015-63502-C3-1-R) (MINECO/FEDER) funded by Fondo Europeo de Desarrollo Regional and Ministerio de Econom/ y Competitividad

References

- Abrahão SM, Gravino C, Pelozo EI, Scanniello G, Tortora G (2013) Assessing the effectiveness of sequence diagrams in the comprehension of functional requirements: Results from a family of five experiments. *IEEE Trans Softw Eng* 39(3):327–342
- Agarwal R, Sinha AP (2003) Object-oriented modeling with UML: a study of developers' perceptions. *Commun ACM* 46(9):248–256
- Anda B, Hansen K, Gullesen I, Thorsen HK (2006) Experiences from introducing UML-based development in a large safety-critical project. *Empir Softw Eng* 11(4):555–581
- Antoniol G, Canfora G, Casazza G, De Lucia A, Merlo E (2002) Recovering traceability links between code and documentation. *IEEE Trans Softw Eng* 28(10):970–983
- Arisholm E, Briand LC, Hove SE, Labiche Y (2006) The impact of UML documentation on software maintenance: An experimental evaluation. *IEEE Trans Softw Eng* 32(6):365–381
- Basili VR, Rombach HD (1988) The TAME project: Towards improvement-oriented software environments. *IEEE Trans Softw Eng* 14(6):758–773
- Basili V, Shull F, Lanubile F (1999) Building knowledge through families of experiments. *IEEE Trans Softw Eng* 25(4):456–473
- Bavota G, Canfora G, Di Penta M, Oliveto R, Panichella S (2013) An empirical investigation on documentation usage patterns in maintenance tasks. In: *Proceedings of International Conference on Software Maintenance*. IEEE Computer Society, pp 210–219
- Bruegge B, Dutoit AH (2003) *Object-oriented software engineering: using UML, Patterns and Java*, 2nd edn. Prentice-Hall, Upper Saddle River
- Budgen D, Burn AJ, Brereton OP, Kitchenham B, Pretorius R (2011) Empirical evidence about the UML: a systematic literature review. *Softw Pract Exper* 41(4):363–392
- Cariou E, Beugnard A, Jezequel JM (2002) An architecture and a process for implementing distributed collaborations. In: *Proceedings of International Enterprise Distributed Object Computing*, pp 132–143
- Carver J, Jaccheri L, Morasca S, Shull F (2003) Issues in using students in empirical studies in software engineering education. In: *Proceedings of International Symposium on Software Metrics*. IEEE Computer Society, pp 239–250
- Corazza A, Maggio V, Scanniello G (2016) Coherence of comments and method implementations: a dataset and an empirical investigation. *Softw Q J*:1–27. <https://doi.org/10.1007/s11219-016-9347-1>
- Dzidek WJ, Arisholm E, Briand LC (2008) A realistic empirical evaluation of the costs and benefits of UML in software maintenance. *IEEE Trans Softw Eng* 34(3):407–432
- Eclipse Modeling Framework (EMF) (2012) <http://www.eclipse.org/modeling/emf/>
- Erickson J, Siau K (2007) Theoretical and practical complexity of modeling methods. *Commun ACM* 50(8):46–51
- Fernández-Sáez AM, Chaudron MRV, Genero M (2013) Exploring costs and benefits of using UML on maintenance: Preliminary findings of a case study in a large IT department. In: *Proceedings of the International Workshop on Experiences and Empirical Studies in Software Modeling co-located with the International Conference on Model Driven Engineering Languages and Systems*, pp 33–42
- Fernández-Sáez AM, Caivano D, Genero M, Chaudron MRV (2015) On the use of UML documentation in software maintenance: Results from a survey in industry. In: *Proceedings of ACM/IEEE International Conference on Model Driven Engineering Languages and Systems*, pp 292–301
- Fernández-Sáez AM, Genero M, Caivano D, Chaudron MRV (2016) Does the level of detail of UML diagrams affect the maintainability of source code?: a family of experiments. *Empir Softw Eng* 21(1):212–259
- Fernández-Sáez AM, Genero M, Chaudron MRV (2013) Empirical studies concerning the maintenance of UML diagrams and their use in the maintenance of code: A systematic mapping study. *Inf Softw Technol* 55(7):1119–1142
- Fernández-Sáez AM, Genero M, Chaudron MRV, Caivano D, Ramos I (2015) Are forward designed or reverse-engineered UML diagrams more helpful for code maintenance?: A family of experiments. *Inf Softw Technol* 57:644–663
- Fluri B, Wursch M, Gall H (2007) Do code and comments co-evolve? on the relation between source code and comment changes. In: *Proceedings of the Working Conference on Reverse Engineering*. IEEE Computer Society, pp 70–79
- Fu R, Gartlehner G, Grant M, Shamliyan T, Sedrakyan A, Wilt TJ, Griffith L, Oremus M, Raina P, Ismaila A, Santaguida P, Lau J, Trikalinos TA (2011) Conducting quantitative synthesis when comparing medical interventions: AHRQ and the effective health care program. *J Clin Epidemiol* 64(11):1187–1197
- Gamma E, Helm R, Johnson R, Vlissides J (1995) *Design patterns: elements of reusable object oriented software*. Addison-Wesley, Boston

- Garousi G, Garousi V, Moussavi M, Ruhe G, Smith B (2013) Evaluating usage and quality of technical software documentation: an empirical study. In: Proceedings of International Conference on Evaluation and Assessment in Software Engineering, pp 24–35
- Gravino C, Tortora G, Scanniello G (2010) An empirical investigation on the relation between analysis models and source code comprehension. In: Proceedings of the International Symposium on Applied Computing. ACM, pp 2365–2366
- Gravino C, Scanniello G, Tortora G (2015) Source-code comprehension tasks supported by UML design models: Results from a controlled experiment and a differentiated replication. *J Vis Lang Comput* 28:23–38
- Grossman M, Aronson JE, McCarthy RV (2005) Does UML make the grade? Insights from the software development community. *Inf Softw Technol* 47(6):383–397
- Guéhéneuc YG (2007) P-mart: Pattern-like micro architecture repository. In: Proceedings of EuroPLoP Focus Group on Pattern Repositories
- Hammad M, Collard ML, Maletic JI (2011) Automatically identifying changes that impact code-to-design traceability during evolution. *Softw Qual J* 19(1):35–64
- Higgins JPT, Green S (2008) *Cochrane handbook for systematic reviews of interventions*, 5th edn. The Cochrane Collaboration, London
- Huedo-Medina TB, Sánchez-Meca J, Marín-Martínez F, Botella J (2006) Assessing heterogeneity in meta-analysis: Q statistic or I^2 index? *Psychol Methods* 11(2):193–206
- Hutchinson JE, Whittle J, Rouncefield M, Kristoffersen S (2011) Empirical assessment of MDE in industry. In: Proceedings of the International Conference on Software Engineering, pp 471–480
- ISO (1991) *Information Technology—Software Product Evaluation: Quality Characteristics and Guidelines for their Use ISO/IEC IS 9126*. ISO, Geneva
- ISO (2000) *ISO 9241-11: Ergonomic requirements for office work with visual display terminals (VDTs) – Part 9: Requirements for non-keyboard input devices*. ISO, Geneva
- ISO (2011) *ISO/IEC 25010 Systems and software engineering – Systems and software Quality Requirements and Evaluation (SQuaRE) – System and software quality models*. ISO, Geneva
- Jedlitschka A, Ciolkowski M, Pfahl D, Sjöberg D (2008) Reporting experiments in software engineering. In: Shull F, Singer J (eds) *Guide to Advanced Empirical Software Engineering*. Springer, pp 201–228
- Jiang ZM, Hassan AE (2006) Examining the evolution of code comments in postgresql. In: Proceedings of Mining Software Repositories. ACM, pp 179–180
- Juristo N, Moreno A (2001) *Basics of software engineering experimentation*. Kluwer Academic Publishers, Dordrecht
- Kitchenham B, Pfleeger S, Pickard L, Jones P, Hoaglin D, El Emam K, Rosenberg J (2002) Preliminary guidelines for empirical research in software engineering. *IEEE Trans Softw Eng* 28(8):721–734
- Kitchenham B, Charters S (2007) Guidelines for performing systematic literature reviews in software engineering
- Lehnert S, Farooq Qua, Riebisch M (2013) Rule-based impact analysis for heterogeneous software artifacts. In: Proceedings of the European Conference on Software Maintenance and Reengineering, pp 209–218
- Leotta M, Ricca F, Antoniol G, Garousi V, Zhi J, Ruhe G (2013) A pilot experiment to quantify the effect of documentation accuracy on maintenance tasks. In: Proceedings of International Conference on Software Maintenance, pp 428–431
- OMG (2005) *Unified modeling language (UML) specification version 2.0*. Technical report, Object Management Group
- Oppenheim AN (1992) *Questionnaire design, interviewing and attitude measurement*. Pinter, London
- Pavalkis S, Nemuraite L (2013) Process for applying derived property based traceability framework in software and systems development life cycle. Springer, Berlin Heidelberg, pp 122–133
- Pavalkis S, Nemuraite L, Butkiene R (2013) Derived properties: A user friendly approach to improving model traceability. *Inf Technol Control* 42(1):48–60
- Pickard L, Kitchenham B, Jones P (1998) Combining empirical results in software engineering. *Inf Softw Technol* 40(14):811–821
- Ried K (2008) *Interpreting and understanding meta-analysis graphs - A practical guide*, vol 35. Australian College of General Practitioners, East Melbourne
- Salviulo F, Scanniello G (2014) Dealing with identifiers and comments in source code comprehension and maintenance: Results from an ethnographically-informed study with students and professionals. In: Proceedings of International Conference on Evaluation and Assessment in Software Engineering. ACM
- Scanniello G, Gravino C, Risi M, Tortora G (2010) A controlled experiment for assessing the contribution of design pattern documentation on software maintenance. In: Proceedings of the Symposium on Empirical Software Engineering and Measurement. ACM

- Scanniello G, Gravino C, Tortora G (2010) Investigating the role of UML in the software modeling and maintenance - a preliminary industrial survey. In: Proceedings of International Conference on Enterprise Information Systems, pp 141–148
- Scanniello G, Gravino C, Genero M, Cruz-Lemus JA, Tortora G (2014) On the impact of UML analysis models on source code comprehensibility and modifiability. *ACM Trans Softw Eng Methodol* 23(2):13:1–13:26
- Scanniello G, Gravino C, Risi M, Tortora G, Dodero G (2015) Documenting design-pattern instances: A family of experiments on source-code comprehensibility. *ACM Trans Softw Eng Methodol* 24(3):14
- Scanniello G, Gravino C, Tortora G, Genero M, Risi M, Cruz-Lemus JA, Dodero G (2015) Studying the effect of uml-based models on source-code comprehensibility: Results from a long-term investigation. In: Springer (ed.) Proceedings of International Conference on Product-Focused Software Process Improvement, vol 9459. Lecture Notes in Computer Science, pp 311–327
- Settimi R, Cleland-Huang J, Khadra OB, Mody J, Lukasik W, DePalma C (2004) Supporting software evolution through dynamically retrieving traces to uml artifacts. In: Proceedings of International Workshop on Principles of Software Evolution, pp 49–54
- Shull F, Carver JC, Vegas S, Juzgado NJ (2008) The role of replications in empirical software engineering. *Empir Softw Eng* 13(2):211–218
- Sillito J, Murphy GC, De Volder K (2008) Asking and answering questions during a programming change task. *IEEE Trans Softw Eng* 34(4):434–451
- Tang A, Jin Y, Han J (2007) A rationale-based architecture model for design traceability and reasoning. *J Syst Softw* 80(6):918–934
- Tang A, Nicholson A, Jin Y, Han J (2007) Using bayesian belief networks for change impact analysis in architecture design. *J Syst Softw* 80(1):127–148
- Wohlin C, Runeson P, Höst M., Ohlsson M, Regnell B, Wesslén A (2012) Experimentation in software engineering. Springer, Berlin
- Zhi J, Sun B, Garousi G, Shahnewaz SM, Ruhe G (2015) Cost, benefits and quality of software development documentation: A systematic mapping. *J Syst Softw* 99:175–198



Giuseppe Scanniello received his Laurea and Ph.D. degrees, both in Computer Science, from the University of Salerno, Italy, in 2001 and 2003, respectively. In 2006, he joined, as an Assistant Professor, the Department of Mathematics and Computer Science at the University of Basilicata, Potenza, Italy. In 2015, he became an Associate Professor at the same university. His research interests include requirements engineering, empirical software engineering, reverse engineering, reengineering, software visualization, workflow automation, migration, wrapping, integration, testing, green software engineering, global software engineering, cooperative supports for software engineering, visual languages and e-learning. He has published more than 150 referred papers in journals, books, and conference proceedings. He serves on the organizing of major international conferences (as general chair, program co-chair, proceedings chair, and member of the program committee) and workshops in the field of software engineering (e.g., ICSE, ASE, ICSME, ICPC, SANER, and many others). Giuseppe Scanniello leads both the group and the laboratory of software engineering at the University of Basilicata (BASELab). He recently obtained the Italian National Scientific Qualification as Full Professor in Computer Science. He is a member of IEEE and IEEE Computer Society. More on <http://www2.unibas.it/gscanniello/>.



Carmine Gravino is an Associate Professor at the Department of Computer Science of University of Salerno. Since 2006 he has been teaching several courses, including Software Engineering, O-O Programming, and Software Metrics and Quality. He is the co-director of the Software Quality and Measurement (SQM)/Web Engineering Laboratory and his research interests include software project management, software measurement and functional size measurement methods, predictive modelling for software engineering, software maintenance and evolution, software technology evaluation through experimental means. He has published more than 100 papers in international journals, books, and conference proceedings and received several awards. He has served as organizing and program committee member of several international conferences in the field of software engineering and is in the editorial boards of international journals. He has also served as reviewer of several software engineering journals. He is a member of the IEEE.



Marcela Genero is Full Professor at the Department of Technologies and Information Systems at the University of Castilla-La Mancha, Ciudad Real, Spain. She received her MSc degree in Computer Science in the Department of Computer Science of the University of South, Argentine in 1989, and her PhD at the University of Castilla-La Mancha, Ciudad Real, Spain in 2002. Her research focuses on the following areas: empirical software engineering, software quality, quality models, conceptual models quality, big data quality, software modelling effectiveness, etc.. Marcela Genero has published more than 100 peer review papers in prestigious journals (DKE, EMSE, ACM TOSEM, IST, JSS, SOSYM, etc.) and conferences (CAISE, E/R, MODELS, ISESE, METRICS, ESEM, EASE, etc.). She co-edited the books titled “Data and Information Quality” (Kluwer, 2001) and “Metrics for Software Conceptual Models” (Imperial College, 2005), among others. She participated in several program committees (CAISE, EASE, ESEM, ICEIS, ICSM, ISESE, etc.) and as reviewer of several journals as well (DKE, EMSE, IEEE TSE, JSS, IST, SOSYM, etc.). She has organised several conferences, workshops and tutorials on empirical studies in software modelling, evidence-based software engineering, quality in conceptual modelling, etc. She has managed several research projects which involved universities and private companies as partners, related to topics within the research areas previously mentioned. She is member of the International Software Engineering Research Network (ISERN) since 2004.



José A. Cruz-Lemus is an Associate Professor at the Department of Information Systems and Technologies at the University of Castilla-La Mancha, Ciudad Real, Spain. He is PhD in Computer Science from the same university. His main research interests are empirical software engineering, software metrics and UML models quality. He has published his works in several journals (Empirical Software Engineering, Information Sciences, Information and Software Technologies, Journal of Systems and Software, etc.), conferences (MoDELS, E/R, ESEM, etc.), and several book chapters.



Genoveffa Tortora received the Laurea degree in Computer Science from the University of Salerno. From 1978 to 1998, she has been with the Department of Computer Scienze. Since 1990, she is a full professor of computer science. In 1998 she founded the Department of Mathematics and Computer Science and has been department chair. From 2000 to 2008 she has been Dean of the Faculty of Mathematical, Physical and Natural Sciences of the University of Salerno. She is an associate editor of International Journal of Software Engineering and Knowledge Engineering, and associate editor of Journal of Visual Languages and Computing. From November 1999 to October 2000, she was a member of the board of directors (Consiglio di Amministrazione) of the University of Salerno. From November 2015 she is a member of the board of directors (Consiglio di Amministrazione) of Fondazione Ravello. From December 2015 she is a member of CNGR (National Committee of Research Guarantors) of the Italian Ministry of Education, University and Research. Her research interests include software engineering, image processing and biometric systems, human-computer interaction, visual languages, databases, datawarehouses, geographic information systems.



Michele Risi received the Laurea degree in computer science in 2001 and the PhD degree in computer science from the University of Salerno, Italy, in 2005. He is Assistant Professor in the Department of Computer Science at the University of Salerno since 2016. His research interests include Reverse Engineering (architecture and design pattern recovery), Reengineering, Human-Computer Interaction, Empirical Software Engineering, Big-Data Analysis, Data-warehouse and Data Visualization, Visual Languages (visual programming environment, parsing techniques and sketch understanding), and Mobile Development and Applications. He has published 90 papers on these topics in international journals, books, and conference and workshop proceedings. He has served as program committee member of several international conferences, and he is member of the Review Board and Editorial Board of international journals.



Gabriella Dodero retired in 2017, was Full Professor and Vice Rector for Studies at the Free University of Bozen-Bolzano. Prior to joining FUB in 2006, she was Associate Professor of Computer Science at the University of Genova, Italy, where she obtained a degree in Mathematics in 1977. Her research interests deal with empirical software engineering, technologies for education and open educational resources. In her career, Prof Dodero has also served in various offices at both universities, as Dean of the Faculty of Computer Science (2007-2009), as Director of University/Enterprise Training Partnership at the University of Genova (1992-1996), and as Vice President of ISICT, a regional consortium for excellency in higher education sponsoring top students in ICT subjects (2003-2006).