

A Controlled Experiment for Program Comprehension through Trace Visualization

Bas Cornelissen, Andy Zaidman, and Arie van Deursen

Report TUD-SERG-2009-001

TUD-SERG-2009-001

Published, produced and distributed by:

Software Engineering Research Group
Department of Software Technology
Faculty of Electrical Engineering, Mathematics and Computer Science
Delft University of Technology
Mekelweg 4
2628 CD Delft
The Netherlands

ISSN 1872-5392

Software Engineering Research Group Technical Reports:

<http://www.se.ewi.tudelft.nl/techreports/>

For more information about the Software Engineering Research Group:

<http://www.se.ewi.tudelft.nl/>

Note: Accepted for publication in IEEE Transactions on Software Engineering.

© copyright 2009, by the authors of this report. Software Engineering Research Group, Department of Software Technology, Faculty of Electrical Engineering, Mathematics and Computer Science, Delft University of Technology. All rights reserved. No part of this series may be reproduced in any form or by any means without prior written permission of the authors.

A Controlled Experiment for Program Comprehension through Trace Visualization

Bas Cornelissen, Andy Zaidman, *Member, IEEE Computer Society*,
and Arie van Deursen, *Member, IEEE Computer Society*

Abstract

Software maintenance activities require a sufficient level of understanding of the software at hand that unfortunately is not always readily available. Execution trace visualization is a common approach in gaining this understanding, and among our own efforts in this context is EXTRA-VIS, a tool for the visualization of large traces. While many such tools have been evaluated through case studies, there have been no quantitative evaluations to the present day. This paper reports on the first controlled experiment to quantitatively measure the added value of trace visualization for program comprehension. We designed eight typical tasks aimed at gaining an understanding of a representative subject system, and measured how a control group (using the Eclipse IDE) and an experimental group (using both Eclipse and EXTRA-VIS) performed in terms of correctness and time spent. The results are statistically significant in both regards, showing a 22% decrease in time needed for the given tasks, and a 43% increase in correctness of the results for the group using trace visualization.

Index Terms

Program comprehension, dynamic analysis, controlled experiment.

B. Cornelissen is with the Software Improvement Group, A.J. Ernststraat 595-H, 1082 LD Amsterdam, The Netherlands.
E-mail: b.cornelissen@sig.eu.

A. Zaidman and A. van Deursen are with the Faculty of Electrical Engineering, Mathematics and Computer Science, Delft University of Technology, Mekelweg 4, 2628CD Delft, The Netherlands.
E-mail: {a.e.zaidman, arie.vandeursen}@tudelft.nl.

I. INTRODUCTION

Program comprehension has become an increasingly important aspect of the software development process. As software systems grow larger and their development becomes more expensive, they are constantly modified rather than built from scratch, which means that a great deal of effort is spent on performing maintenance activities. However, as up to date documentation is often lacking, it is estimated that up to 60% of the maintenance effort is spent on gaining a sufficient *understanding* of the program at hand [1], [2]. It is for this reason that the development of techniques and tools that support the comprehension process can make a significant contribution to the overall efficiency of software development.

With respect to such techniques, the literature offers numerous solutions that can be roughly broken down into static and dynamic approaches (and combinations thereof). Whereas static analysis relies on such artifacts as source code and documentation, dynamic analysis focuses on a system's execution. An important advantage of dynamic analysis is its precision, as it captures the system's actual behavior. Among the drawbacks are its incompleteness, as the gathered data pertains solely to the scenario that was executed; and the well-known scalability issues, due to the often excessive amounts of execution trace data. Particularly this latter aspect is troublesome because of the cognitive overload on the part of the maintainer.

To cope with the issue of scalability, a significant portion of the literature on program comprehension has been dedicated to the reduction [3], [4] and visualization [5], [6] of execution traces. One of these techniques and tools is EXTRAVIS, a tool that offers two interactive views of large execution traces [7]. Through a series of case studies we illustrated how EXTRAVIS can support different types of common program comprehension activities. However, in spite of these efforts, there is no *quantitative* evidence of the tool's usefulness in practice. As we will show in the next section, no such evidence is offered for any of the trace visualization techniques in the program comprehension literature.

The purpose of this paper, therefore, is a first quantification of the usefulness of trace visualization for program comprehension. Furthermore, to gain a deeper understanding of the nature of its added value, we investigate which types of tasks benefit most from trace visualization and from EXTRAVIS. To fulfill these goals, we design and execute a controlled experiment in which we measure how the tool affects (1) the time that is needed for typical comprehension tasks,

and (2) the correctness of the solutions given during those tasks.

This paper extends our previous work [8] with a survey of 21 trace visualization techniques, an additional group of subjects with an industrial background (thus strengthening the statistical significance as well as the external validity), and a discussion on the implications of our EXTRAVIS findings for trace visualization tools in general.

The remainder of the paper is structured as follows. Section II extensively reviews existing techniques and tools for trace visualization, and motivates our intent to conduct a controlled experiment. Section III offers a detailed description of the experimental design. Section IV presents the results of our experiment, which are then discussed in Section V. Section VI discusses threats to validity, and Section VII offers conclusions and future directions.

II. BACKGROUND

A. Execution trace analysis

The use of dynamic analysis for program comprehension has been a popular research activity in the last decades. In a large survey that we recently performed, we identified a total of 176 articles on this topic that were published between 1972 and June 2008 [9]. More than 30 of these papers concern *execution trace analysis*, which has often shown to be beneficial to such activities as feature location, behavioral analysis, and architecture recovery.

Understanding a program through its execution traces is not an easy task because traces are typically too large to be comprehended directly. Reiss and Renieris, for example, report on an experiment in which one gigabyte of trace data was generated for every two seconds of executed C/C+ code or every ten seconds of Java code [3]. For this reason, there has been a significant effort in the automatic *reduction* of traces to make them more tractable (e.g., [3], [10], [4]). The reduced traces can then be *visualized* by traditional means: for example, as directed graphs or UML sequence diagrams. On the other hand, the literature also offers several non-traditional trace visualizations that have been designed specifically to address the scalability issues.

In Section II-B we present an overview of the current state of the art in trace visualization. Section II-C describes EXTRAVIS, our own solution, and Section II-D motivates the need for controlled experiments.

TABLE I
OVERVIEW OF EXISTING TRACE VISUALIZATION TECHNIQUES

References	Tool	Evaluation type	Applications
[11]	GRAPHTRACE	small case study	debugging
[5], [12], [13], [14]	JINSIGHT; OVATION; TPTP*	preliminary; user feedback	general understanding
[15]	SCENE*	preliminary	software reuse
[6], [16]	ISVIS*	case study	architecture reconstruction, feature location
[17], [18]	SCED; SHIMBA	case study	debugging; various comprehension tasks
[19]	FORM	case study	detailed understanding; distributed systems
[20]	JAVAVIS	preliminary; user feedback	educational purposes; detailed understanding
[21], [4], [22], [23]	SEAT	small case studies; user feedback	general understanding
[24], [25], [26], [27]	SCENARIOGRAPHER	multiple case studies	detailed understanding; distributed systems; feature analysis; large-scale software
[28], [29], [30]	–	small case study	quality control; conformance checking
[10]	–	multiple case studies	general understanding
[31]	–	case study	trace comparison; feature analysis
[32]	–	case study	feature analysis
[33]	–	case study	architecture reconstruction; conformance checking; behavioral profiles
[34]	TRACEGRAPH	industrial case study	feature analysis
[35], [36]	SDR; JRET*	multiple case studies	detailed understanding through test cases
[37]	FIELD; JIVE; JOVE	multiple case studies	performance monitoring; phase detection
[38]	–	–	API understanding
[39], [7]	EXTRAVIS*	multiple case studies	fan-in/-out analysis; feature analysis; phase detection
[40]	OASIS	user study	various comprehension tasks
[41]	–	small case studies	general understanding; wireless sensor networks

B. Execution trace visualization

There exist three surveys in the area of execution trace visualization that provide overviews of existing techniques. The first survey was published in 2003 by Pacione et al., who compare the performance of five dynamic visualization tools [42]. Another survey was published in 2004 by Hamou-Lhadj and Lethbridge, who describe eight trace visualization tools from the literature [43]. Unfortunately, these two overviews are incomplete because (1) the selection procedures were non-systematic, which means that papers may have been missed; and (2) many more solutions have been proposed in the past five years. A third survey was performed by the authors of this paper in 2008, and was set up as a large-scale systematic literature survey of all dynamic analysis-based approaches for program comprehension [9]. However, its broad perspective prevents subtle differences between trace visualization techniques from being exposed, particularly in terms of evaluation: for example, it does not distinguish between user studies and controlled experiments.

To obtain a complete overview of all existing techniques and to reveal the differences in evaluation, we have used our earlier survey to identify all articles on trace visualization for program comprehension from 1988 onwards, and then reexamined these papers from an evaluation perspective. In particular, we have focused on techniques that *visualize (parts of) execution traces*. We have looked at the types of validation and the areas in which the techniques were applied. Also of interest was the public availability of the tools involved, which is crucial for fellow researchers seeking to study existing solutions or perform replications of the experiment described in this paper.

Our study has resulted in the identification and characterization of 21 contributions¹ that were published between 1988 and 2008, shown in Table I. For each contribution, the table shows the appropriate references, associated tools (with asterisks denoting public availability), evaluation types, and areas in which the technique was applied. In what follows, we briefly describe the contents of each paper.

1988-2000

Kleyn and Gingrich were among the first to point out the value of visualizing run-time behavior [11]. Their visualization of execution traces is graph-based and aims at better understanding

¹Of the 36 papers found, Table I shows only the 21 unique contributions (i.e., one per first author).

software and identifying programming errors. In particular, their graph visualization is animated, in the sense that the user of the tool can step through the entire execution and observe what part(s) of the program are currently active. A case study illustrates how their views can provide more insight into the inner workings of a system.

De Pauw et al. introduced their interaction diagrams, very similar to UML sequence diagrams, in Jinsight, a tool which could visualize running Java programs [5]. Jinsight was later transformed into the publicly available TPTP Eclipse plugin, which brings execution trace visualization to the mainstream Java developer. The authors also noticed that the standard sequence diagram notation was difficult to scale up for large software systems, leading to the development of their “*execution pattern*” notation, a much more condensed view of the typical sequence diagram [12].

Koskimies and Mössenböck proposed Scene, which combines a sequence diagram-like visualization with hypertext-like facilities [15]. These hypertext facilities allow the user to browse related documents such as the source code or UML class diagrams. The authors are aware of scalability issues when working with sequence diagrams and therefore proposed a number of abstraction techniques.

Jerding et al. created the “*Interaction Scenario Visualizer*” (ISVis) [6], [16]. ISVis combines static and dynamic information to accomplish amongst others *feature location*, the establishment of relations between concepts and source code [44]. ISVis’ dynamic component visualizes *scenario views*, which bear some resemblance to sequence diagrams. Of particular interest is the *Information Mural* view, which effectively provides an overview of an entire execution scenario, comprising hundreds of thousands of interactions. The authors have applied ISVis to the Mosaic web browser in an attempt to extend it.

Systä et al. presented an integrated reverse engineering environment for Java that uses both static and dynamic analysis [17], [18]. The dynamic analysis component of this environment, SCED, visualizes the execution trace as a sequence diagram. In order to validate their approach, a case study was performed on the Fujaba open source UML tool suite, in which a series of program comprehension and reverse engineering tasks were conducted.

2000-2005

Souder et al. were among the first to recognize the importance of understanding distributed

applications with the help of dynamic analysis [19]. To this purpose, they use Form, which enables to draw sequence diagrams for distributed systems. The authors validate their approach through a case study.

Oeschle and Schmitt built a tool called JAVAVIS that visualizes running Java software, amongst others through sequence diagrams [20]. The authors' main aim was to use JAVAVIS for educational purposes and their validation comprises informal feedback from students using the tool.

Hamou-Lhadj et al. created the Software Exploration and Analysis Tool (SEAT) that visualizes execution traces as trees. It is integrated in the IDE to enable easy navigation between different views [22]. SEAT should be considered as a research vehicle in which the authors explored some critical features of trace visualization tools. Subsequently, they began exploring such solutions, such as trace compression [4] or removing parts of the trace without affecting its overall information value [23]. While the degree of compression is measured in several case studies, the added value for program comprehension remains unquantified.

Salah and Mancoridis investigate an environment that supports the comprehension of distributed systems, which are typically characterized by the use of multiple programming languages [24]. Their environment visualizes sequence diagrams, with a specific notation for inter-process communication. The authors also report on a small case study. Salah et al. later continued their dynamic analysis work and created the so-called module-interaction view, that shows which modules are involved in the execution of a particular use case [27]. They evaluate their visualization in a case study on Mozilla and report on how their technique enables feature location.

Briand et al. specifically focused on visualizing sequence diagrams from distributed applications [28], [30]. Through a small case study with their prototype tool they have reverse engineered sequence diagrams for checking design conformance, quality, and implementation choices.

Zaidman and Demeyer represented traces as signals in time [10]. More specifically, they count how many times individual methods are executed and using this metric, they visualize the execution of a system throughout time. This allows to identify phases and re-occurring behavior. They show the benefits of their approach using two case studies.

2006-2007

Kuhn and Greevy also represented traces as signals in time with their "dynamic time warping"

approach [31]. In contrast to Zaidman and Demeyer, they rely on the stack depth as the underlying metric. The signals are compared to one another to locate features, as illustrated by a case study.

Greevy et al. explored polymetric views to visualize the behavior of features [32]. Their 3D visualization renders run-time events of a feature as towers of instances, in which a tower represents a class and the number of boxes that compose the tower indicates the number of live instances. Message sends between instances are depicted as connectors between the boxes. The authors perform a case study to test their approach.

Koskinen et al. proposed *behavioral profiles* to understand and identify extension points for components [33]. Their technique combines information from execution traces and behavioral rules defined in documentation to generate these profiles, which contain an architectural level view on the behavior of a component or application. Their ideas are illustrated in a case study.

Simmons et al. used TraceGraph to compare execution traces with the aim of locating features [34]. Furthermore, they integrate the results of their feature location technique into a commercial static analysis tool so as to make feature location more accessible to their industrial partner. The authors furthermore report on a case study performed in an industrial context.

2007-2008

Cornelissen et al. looked specifically into generating sequence diagrams from test cases, arguing that test scenarios are relatively concise execution scenarios that reveal a great deal about the system's inner workings [35]. They initially applied their SDR tool to a small case study, and later extended their ideas in the publicly available JRET eclipse plugin, which was evaluated on a medium-scale open source application [36].

Over the years, Reiss has developed numerous solutions for visualizing run-time behavior [37]. Among the most notable examples are FIELD, which visualizes dynamic call graphs, and JIVE, which visualizes the execution behavior in terms of classes or packages. JIVE's visualization breaks up time in intervals and for each interval it portrays information such as the number of allocations, the number of calls, and so on.

Jiang et al. concentrated on generating sequence diagrams specifically for studying API usage [38]. The rationale of their approach is that it is often difficult to understand how APIs should be used or can be reused. An evaluation of their approach is as yet not available.

Bennett et al. engineered the Oasis Sequence Explorer [40]. Oasis was created based on a

focus group experiment that highlighted some of the most desirable features when exploring execution traces. The authors then performed a user study to validate whether the Oasis features were indeed helpful during a series of typical software maintenance tasks, with quite useful measurements as a result.

Dalton and Hallstrom designed a dynamic analysis visualization toolkit specifically aimed at TinyOS, a component-based operating system mainly used in the realm of wireless sensor networks [41]. They generate annotated call graphs and UML sequence diagrams for studying and understanding TinyOS applications. They illustrate the benefits of their tool through a case study on a TinyOS component.

C. *Extravis*

Among our own contributions to the field of trace visualization is EXTRAVIS. This publicly available² tool provides two linked, interactive views, shown in Figure 1. The *massive sequence view* is essentially a large-scale UML sequence diagram (similar to Jerding's Information Mural [45]), and offers an overview of the trace and the means to navigate it. The *circular bundle view* hierarchically projects the program's structural entities on a circle and shows their inter-relationships in a bundled fashion. A comparison of EXTRAVIS with other tools is provided in our earlier work [7].

We qualitatively evaluated the tool in various program comprehension contexts, including trace exploration, feature location, and top-down program comprehension [7]. The results provided initial evidence of EXTRAVIS' benefits in these contexts, the main probable advantages being its optimal use of screen real estate and the improved insight into a program's structure. However, we hypothesized that the relationships in the circular view may be difficult to grasp.

D. *Validating trace visualizations*

The overview in Table I shows that trace visualization techniques in the literature have been almost exclusively evaluated using case studies. Indeed, there have been no efforts to *quantitatively* measure the usefulness of trace visualization techniques in practice, e.g., through controlled experiments. Moreover, the evaluations in existing work rarely involve broad spectra

²EXTRAVIS, <http://swerl.tudelft.nl/extravis>

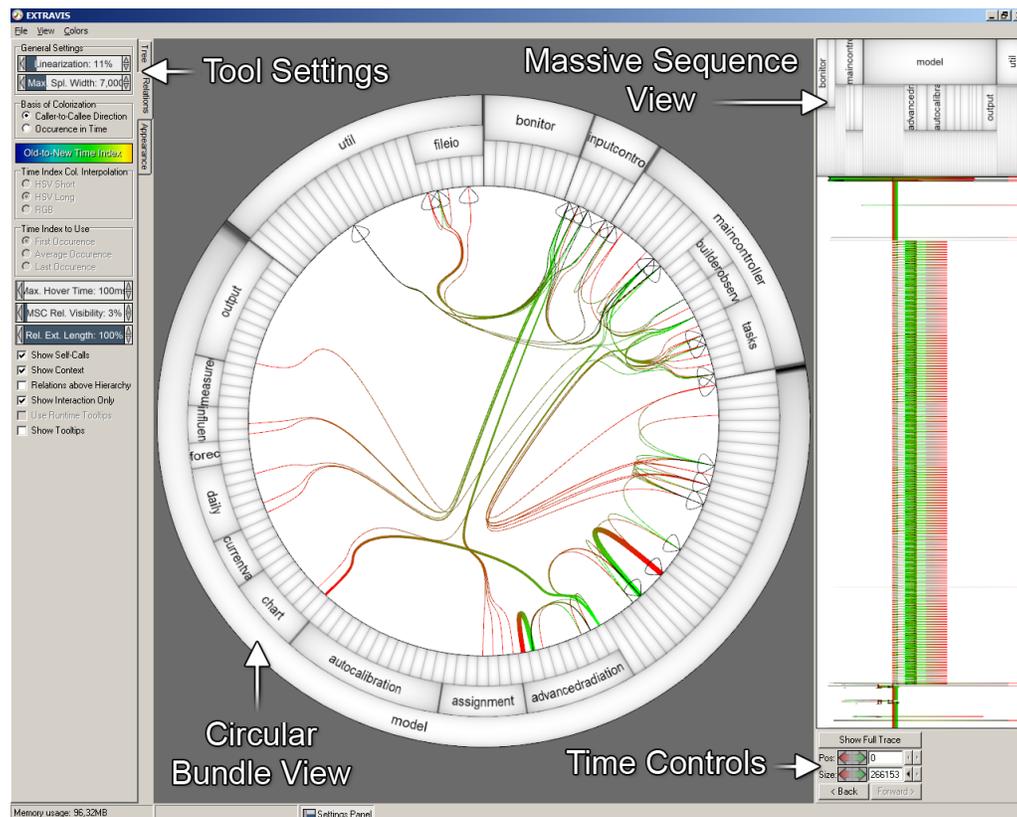


Fig. 1. EXTRAVIS' circular bundle view and massive sequence view.

of comprehension tasks, making it difficult to judge whether the associated solutions are widely applicable in daily practice. Lastly, most existing approaches involve *traditional* visualizations, i.e., they rely on UML, graph, or tree notations, to which presumably most software engineers are accustomed [9]. By contrast, EXTRAVIS uses non-traditional visualization techniques, and Storey argues [46] that advanced visual interfaces are not often used in development environments because they tend to require complex user interactions.

These reasons have motivated us to empirically validate EXTRAVIS through a controlled experiment in which we seek to assess its added value in concrete maintenance contexts.

III. EXPERIMENTAL DESIGN

The purpose of this paper is to provide a quantitative evaluation of trace visualization for program comprehension. To this end, we define a series of typical comprehension tasks and measure EXTRAVIS' added value to a traditional programming environment: in this case, the Eclipse IDE³. Similar to related efforts (e.g., [47], [48]) we maintain a distinction between *time spent* and *correctness*. Furthermore, we seek to identify the types of tasks to which the use of EXTRAVIS, and trace visualization in general, is the most beneficial.

A. Research Questions and Hypotheses

Based on our earlier case studies, we distinguish the following research questions:

- 1) Does the availability of EXTRAVIS reduce the *time* that is needed to complete typical comprehension tasks?
- 2) Does the availability of EXTRAVIS increase the *correctness* of the solutions given during those tasks?
- 3) Based on the answers to these research questions, which *types* of tasks can we identify that benefit most from the use of EXTRAVIS and from trace visualization in general?

Associated with the first two research questions are two null hypotheses, which we formulate as follows:

- $H1_0$: The availability of EXTRAVIS does not impact the time needed to complete typical comprehension tasks.
- $H2_0$: The availability of EXTRAVIS does not impact the correctness of solutions given during those tasks.

The alternative hypotheses that we use in the experiment are the following:

- $H1$: The availability of EXTRAVIS reduces the time needed to complete typical comprehension tasks.
- $H2$: The availability of EXTRAVIS increases the correctness of solutions given during those tasks.

³Eclipse IDE, <http://www.eclipse.org>

The rationale behind the first alternative hypothesis is the fact that EXTRAVIS provides a broad overview of the subject system on one single screen, which may guide the user to his or her goal more quickly than if switching between source files were required.

The second alternative hypothesis is motivated by the inherent precision of dynamic analysis with respect to actual program behavior: for example, the resolution of late binding may result in a more detailed understanding and therefore produce more accurate solutions.

To test hypotheses $H1_0$ and $H2_0$, we define a series of comprehension tasks that are to be addressed by both a control group and an experimental group. The difference in treatment between these groups is that the former group uses a traditional development environment (the “Eclipse” group), whereas the latter group also has access to EXTRAVIS (the “Ecl+Ext” group). We maintain a between-subjects design, meaning that each subject is either in the control or in the experimental group.

Sections III-B through III-G provide a detailed description of the experiment.

B. Object

The system that is to be comprehended by the subject groups is CHECKSTYLE, a tool that employs “checks” to verify if source code adheres to specific coding standards. Our choice for CHECKSTYLE as the object of this experiment is motivated by the following factors:

- CHECKSTYLE is open source, which helps to make the results of our experiments reproducible.
- CHECKSTYLE comprises 310 classes distributed across 21 packages, containing a total of 57 KLOC.⁴ This makes it tractable for an experimental session, yet adequately representative of real life programs.
- It is written in Java, with which many potential subjects are sufficiently familiar.
- It address an application domain, adherence to coding standards, which will be understandable to most potential subjects.
- The authors of this paper are familiar with its internals as a result of earlier experiments [49], [50], [7]. Furthermore, the lead developer is available for feedback.

⁴Measured using `sloccount` by David A. Wheeler, <http://sourceforge.net/projects/sloccount/>.

To obtain the necessary trace data for EXTRAVIS, we instrument CHECKSTYLE and execute it according to two scenarios. Both involve typical runs with a small input source file, and only differ in terms of the input configuration, which in one case specifies 64 types of checks whereas the other specifies only six. The resulting traces contain 31,260 and 17,126 calls, respectively, which makes them too large to be comprehended in limited time without tool support.

Analyzing the cost of creating these traces is not included in the experiment, as it is our prime objective to analyze whether the availability of trace information is beneficial during the program comprehension process. In practice, we suspect that execution traces will likely be obtained from test cases – a route we also explored in our earlier work [35].

C. Task design

With respect to the comprehension tasks that are to be tackled during the experiment, we maintain two important criteria: (1) they should be representative of real maintenance contexts, and (2) they should not be biased towards either Eclipse or EXTRAVIS.

To this end, we apply the comprehension framework from Pacione et al. [51], who argue that “*a set of typical software comprehension tasks should seek to encapsulate the principal activities typically performed during real world software comprehension*”. They have studied several sets of tasks used in software visualization and comprehension evaluation literature and classified them according to nine principal activities, representing both general and specific reverse engineering tasks and covering both static and dynamic information (Table II). Particularly the latter aspect significantly reduces biases towards either of the two tools used in this experiment.

Using these principal activities as a basis, we propose eight representative tasks that highlight many of CHECKSTYLE’s aspects at both high and low abstraction levels. Table III provides descriptions of the tasks and shows how each of the nine activities from Pacione et al. is covered by at least one task.⁵ For example, activity A1, “*Investigating the functionality of (part of) the system*”, is covered by tasks T1, T3.1, T4.1, and T4.2; and activity A4, “*Investigating dependencies between artifacts*”, is covered by tasks T2.1, T2.2, T3.2, and T3.3.

To render the tasks more representative of real maintenance situations, tasks are given as open rather than multiple-choice questions, making it harder for respondents to resort to simply

⁵Table III only contains the actual questions; the subjects were also given contextual information (such as definitions of fan-in and coupling) which can be found in the appendix.

TABLE II
PACIONE'S NINE PRINCIPAL COMPREHENSION ACTIVITIES

Activity	Description
A1	Investigating the functionality of (a part of) the system
A2	Adding to or changing the system's functionality
A3	Investigating the internal structure of an artifact
A4	Investigating dependencies between artifacts
A5	Investigating run-time interactions in the system
A6	Investigating how much an artifact is used
A7	Investigating patterns in the system's execution
A8	Assessing the quality of the system's design
A9	Understanding the domain of the system

guessing. Per answer, 0–4 points can be earned. Points are awarded by the evaluators, in our case the first two authors. A solution model is available (see the appendix), which was reviewed by CHECKSTYLE's lead developer. To ensure uniform grading among the two evaluators, the solution of five random subjects are first graded by both evaluators.

D. Subjects

The subjects in this experiment are fourteen Ph.D. candidates, nine M.Sc. students, three postdocs, two professors, and six participants from industry. The resulting group thus consists of 34 subjects, and is quite heterogeneous in that it represents 8 different nationalities, and M.Sc. degrees from 16 universities. The M.Sc. students are in the final stages of their computer science studies, and the Ph.D. candidates represent different areas of software engineering, ranging from software inspection to model-based fault diagnosis. Our choice of subjects partly mitigates concerns from Di Penta et al., who argue that “*a subject group made up entirely of students might not adequately represent the intended user population*” [52]. All subjects participate on a voluntary basis and can therefore be assumed to be properly motivated. None of them have prior experience with EXTRAVIS.

To partition the subjects, we distinguish five fields of expertise that could strongly influence the individual performance. They represent variables that are to be controlled during the experiment, and concern knowledge of Java, Eclipse, reverse engineering, CHECKSTYLE, and language

TABLE III
DESCRIPTIONS OF THE COMPREHENSION TASKS

Task	Activities	Description
<i>Context: Gaining a general understanding.</i>		
T1	A{1,7,9}	Having glanced through the available information for several minutes, which do you think are the main stages in a typical (non-GUI) Checkstyle scenario? Formulate your answer from a high-level perspective: refrain from using identifier names and stick to a maximum of six main stages.
<i>Context: Identifying refactoring opportunities.</i>		
T2.1	A{4,8}	Name three classes in Checkstyle that have a high fan-in and (almost) no fan-out.
T2.2	A{4,8}	Name a class in the top-level package that could be a candidate for movement to the <code>api</code> package because of its tight coupling with classes therein.
<i>Context: Understanding the checking process.</i>		
T3.1	A{1,2,5,6}	In general terms, describe the life cycle of the <code>checks.whitespace.TabCharacterCheck</code> during execution: when is it created, what does it do and on whose command, and how does it end up?
T3.2	A{3,4,5}	List the identifiers of all method/constructor calls that typically occur between <code>TreeWalker</code> and a <code>TabCharacterCheck</code> instance, and the order in which they are called. Make sure you also take inherited methods/constructors into account.
T3.3	A{3,4,5,9}	In comparison to the calls listed in Task T3.2., which additional calls occur between <code>TreeWalker</code> and <code>checks.coding.IllegalInstantiationCheck</code> ? Can you think of a reason for the difference?
<i>Context: Understanding the violation reporting process.</i>		
T4.1	A{1,3}	How is the check's warning handled, i.e., where/how does it originate, how is it internally represented, and how is it ultimately communicated to the user?
T4.2	A{1,5}	Given <code>Simple.java</code> as the input source and <code>many_checks.xml</code> as the configuration, does <code>checks.whitespace.WhitespaceAfterCheck</code> report warnings? Specify how your answer was obtained.

technology (i.e., CHECKSTYLE's domain). The subjects' levels of expertise in each of these fields are measured through a (subjective) a priori assessment: we use a five-point Likert scale, from 0 ("no knowledge") to 4 ("expert"). In particular, we require minimum scores of 1 for Java and Eclipse ("beginner"), and a maximum score of 3 for CHECKSTYLE ("advanced"). A characterization of the subjects is provided in the appendix.

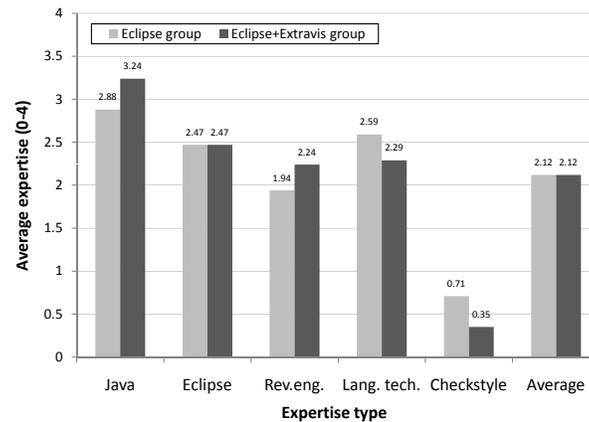


Fig. 2. Average expertise of the subject groups.

The assignments to the control and experimental group are done by hand to evenly distribute the available knowledge. The result is illustrated by Figure 2: in each group, the expertise is chosen to be as similar as possible, resulting in an average expertise of 2.12 in both groups.

E. Experimental procedure

The experiment is performed through a dozen sessions, most of which take place at the university. Sessions with industrial subjects take place at their premises, in our case the Software Improvement Group⁶, the industrial partner in our project. The sessions are conducted on workstations with characteristics that were as similar as possible, i.e., at least Pentium 4 processors and comparable screen resolutions (1280×1024 or 1600×900). Given the different locations (university and in house at company) fully equivalent setups were impossible to achieve.

Each session involves at most three subjects and features a short tutorial on Eclipse, highlighting the most common features. The experimental group is also given a ten minute EXTRAVIS tutorial that involves a JHOTDRAW execution trace used in earlier experiments [7]. All sessions are supervised, enabling the subjects to pose clarification questions, and preventing them from consulting others and from using alternative tools. The subjects are not familiar with the experimental goal.

⁶Software Improvement Group, <http://www.sig.eu>

The subjects are presented with a fully configured Eclipse that is readily usable, and are given access to the example input source file and CHECKSTYLE configurations (see Section III-B). The Ecl+Ext group is also provided with two EXTRAVIS instances, for each of the two execution traces mentioned earlier. All subjects receive handouts that provide an introduction, CHECKSTYLE outputs for the two aforementioned scenarios, the assignment, a debriefing questionnaire, and reference charts for both Eclipse and EXTRAVIS. The assignment is to complete the eight comprehension tasks within 90 minutes. The subjects are required to motivate their answers at all times. We purposely refrain from influencing how exactly the subjects should cope with the time limit: only when a subject exceeds the time limit is he or she told that finishing up is, in fact, allowed. Finally, the questionnaire asks for the subjects' opinions on such aspects as time pressure and task difficulty.

F. Variables & Analysis

The independent variable in our experiment is the availability of EXTRAVIS during the tasks.

The first dependent variable is the *time spent* on each task, and is measured by having the subjects write down the current time when starting a new task. Since going back to earlier tasks is not allowed and the sessions are supervised, the time spent on each task can be easily reconstructed.

The second dependent variable is the *correctness* of the given solutions. This is measured by applying our solution model to the subjects' solutions, which specifies the required elements and the associated scores.

To test our hypotheses, we first test whether the sample distributions are normal (via a Kolmogorov-Smirnov test) and whether they have equal variances (via Levene's test). If these tests pass, we use the parametric Student's t-test to evaluate our hypotheses; otherwise we use the (more robust, but weaker) non-parametric Mann-Whitney test.

Following our alternative hypotheses, we employ the one-tailed variant of each statistical test. For the time as well as the correctness variable we maintain a typical confidence level of 95% ($\alpha=0.05$). The statistical package that we use for our calculations is SPSS.

TABLE IV
DESCRIPTIVE STATISTICS OF THE EXPERIMENTAL RESULTS (34 SUBJECTS VERSION)

	Time		Correctness	
	<i>Eclipse</i>	<i>Ecl+Ext</i>	<i>Eclipse</i>	<i>Ecl+Ext</i>
mean	77.00	59.94	12.47	17.88
difference		-22.16%		+43.38%
min	38	36	5	11
max	102	72	22	22
median	79	66	14	18
stdev.	18.08	12.78	4.54	3.24
one-tailed Student's t-test				
Kolmogorov-Smirnov <i>Z</i>	0.606	0.996	0.665	0.909
Levene <i>F</i>		1.370		2.630
df		32		32
t		3.177		4.000
p-value		0.002		<0.001

G. Pilot studies

Prior to the experimental sessions, we conduct two pilots to optimize several experimental parameters, such as the number of tasks, their clarity, feasibility, and the time limit. The pilot for the control group is performed by an author of this paper who had initially not been involved in the experimental design. The pilot for the experimental group is conducted by an outsider. Both would not take part in the actual experiment later on.

The results of the pilots led to the removal of two tasks because the time limit was too strict. The removed tasks were already taken into account in Section III-B. Furthermore, the studies led to the refinement of several tasks in order to make the questions clearer. Other than these unclarities, the tasks were found to be sufficiently feasible in both the Eclipse and the Ecl+Ext pilot.

IV. RESULTS

Table IV shows descriptive statistics of the measurements, aggregated over all tasks; the measurements themselves are available as a spreadsheet.⁷

⁷Results spreadsheet (34 subjects version), <http://www.st.ewi.tudelft.nl/~cornel/results-new.xlsx>

TABLE V
DESCRIPTIVE STATISTICS OF THE EXPERIMENTAL RESULTS (24 SUBJECTS VERSION)

	Time		Correctness	
	<i>Eclipse</i>	<i>Ecl+Ext</i>	<i>Eclipse</i>	<i>Ecl+Ext</i>
mean	74.75	59.42	12.75	18.25
difference		-20.51%		+43.14%
min	38	36	5	11
max	102	72	19	22
median	78	67	14	19
stdev.	18.34	14.19	4.18	3.25
one-tailed Student's t-test				
Kolmogorov-Smirnov Z	0.512	0.908	0.984	1.049
Levene F		0.467		1.044
df		22		22
t		2.291		3.598
p-value		0.016		0.001
one-tailed Mann-Whitney test				
<i>U</i>		32.50		22.00
p-value		0.011		0.002

Table V shows the same data, but pertaining to our earlier results, involving 24 subjects instead of 34.⁸ The remainder of this report discusses the results of the 34 subjects version.

Wohlin et al. [53] suggest the removal of *outliers* in case of extraordinary situations, such as external events that are unlikely to reoccur. We found four outliers in our timing data and one more in the correctness data, but could identify no such circumstances and have therefore opted to retain those data points.

As an important factor for both time and correctness, we note that two subjects decided to stop after 90 minutes with two tasks remaining, and one subject stopped with one task remaining, resulting in ten missing data points in this experiment (i.e., the time spent by three subjects on task T4.2 and by two subjects on task T4.1, as well as the correctness of the solutions involved). Nine others finished all tasks, but only after the 90 minutes had expired: eight subjects from the Eclipse group and one subject from the Ecl+Ext group spent between 95 and 124 minutes. The

⁸Results spreadsheet (24 subjects version), <http://www.st.ewi.tudelft.nl/~cornel/results.xlsx>

remaining 22 participants finished all eight tasks on time.⁹

In light of the missing data points, we have chosen to disregard the last two tasks in our quantitative analyses. Not taking tasks T4.1 and T4.2 into account, only three out of the 34 subjects still exceeded the time limit (by 6, 7 and 12 minutes, respectively). This approach also has the advantage that any ceiling effects in our data, that may have resulted from the increasing time pressure near the end of the assignment, are strongly reduced. The remaining six tasks still cover all of Pacione's nine activities (Table III).

A. Time results

We start off by testing null hypothesis $H1_0$: the availability of EXTRAVIS does not impact the time that is needed to complete typical comprehension tasks.

Figure 3(a) shows a box plot for the total time that the subjects spent on the first six tasks. Table IV indicates that on average the Ecl+Ext group required 22.16% less time.

The Kolmogorov-Smirnov and Levene tests succeeded for the timing data, which means that Student's t-test may be used to test $H1_0$. As shown in Table IV, the t-test yields a statistically significant result. The average time spent by the Ecl+Ext group was clearly lower and the p-value 0.002 is smaller than 0.05, which means that $H1_0$ can be rejected in favor of the alternative hypothesis $H1$, which states that the availability of EXTRAVIS reduces the time that is needed to complete typical comprehension tasks.

B. Correctness results

We next test null hypothesis $H2_0$, which states that the availability of EXTRAVIS does not impact the correctness of solutions given during typical comprehension tasks.

Figure 3(b) shows a box plot for the scores that were obtained by the subjects on the first six tasks. Note that we consider overall scores rather than scores per task (which are left to Section V-C). The box plot shows that the difference in terms of correctness is even more explicit than for the timing aspect. The solutions given by the Ecl+Ext subjects were 43.38% more accurate (Table IV), averaging 17.88 out of 24 points compared to 12.47 points for the Eclipse group.

⁹Related studies point out that it is not uncommon for several tasks to remain unfinished during the actual experiments (e.g., [48] and [40]).

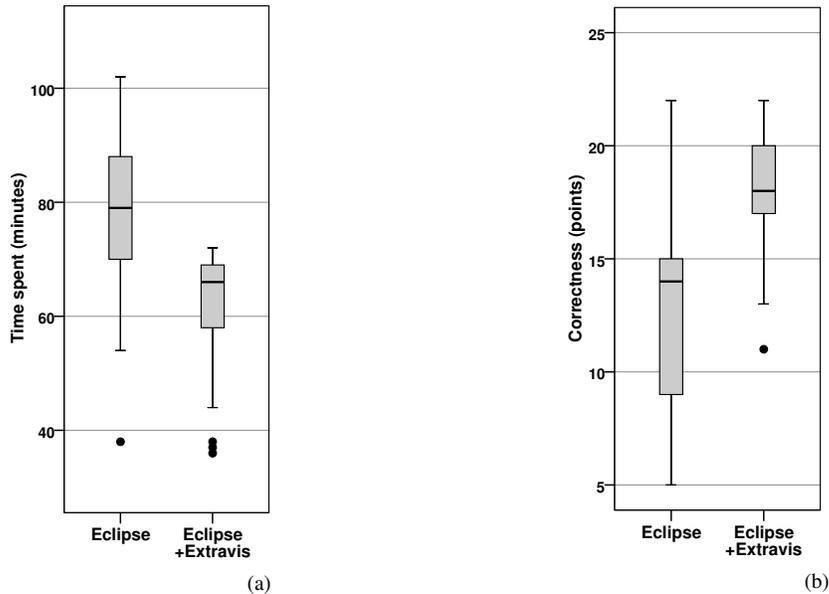


Fig. 3. Box plots for time spent and correctness.

Similar to the timing data, the requirements for the use of the parametric t-test were met. Table IV therefore shows the results for Student's t-test. At less than 0.001, the p-value is very low and implies statistical significance, meaning that H_0 can be rejected in favor of our alternative hypothesis H_2 , which states that the availability of EXTRAVIS increases the correctness of solutions given during typical comprehension tasks.

V. DISCUSSION

A. Reasons for different time requirements

The lower time requirements for the EXTRAVIS users can be attributed to several factors. First, all information offered by EXTRAVIS is shown on a single screen, which eliminates the need for scrolling. In particular, the overview of the entire system's structure saves much time in comparison to conventional environments, in which typically multiple files have to be studied at once. Second, the need to imagine how certain functionalities or interactions work at run-time represents a substantial cognitive load on the part of the user. This is alleviated by trace analysis and visualization tools, which show the actual run-time behavior. Examples of these assumptions will be discussed in Section V-C.

On the other hand, several factors may have had a negative impact on the the time requirements of EXTRAVIS users. For example, the fact that EXTRAVIS is a standalone tool means that context switching is necessary, which may yield a certain amount of overhead on the part of the user. This could be solved by integrating the trace visualization technique into Eclipse (or other IDEs), with the additional benefit that the tool could provide direct links to Eclipse's source code browser. However, it should be noted that EXTRAVIS would still require a substantial amount of screen real estate to be used effectively.

Another potential factor that could have hindered the time performance of the Ecl+Ext group is that these subjects may not have been sufficiently familiar with EXTRAVIS' features, and were therefore faced with a time-consuming learning curve. This is partly supported by the debriefing questionnaire, which indicates that five out of the seventeen subjects found the tutorial too short. A more elaborate tutorial on the use of the tool could help alleviate this issue.

B. Reasons for correctness differences

We attribute the added value of EXTRAVIS to correctness to several factors. A first one is the inherent precision of dynamic analysis: the fact that EXTRAVIS shows the actual objects involved in each call makes the interactions easier to understand. Section V-C discusses this in more detail through an example task.

Second, the results of the debriefing questionnaire (Table VI) show that the Ecl+Ext group used EXTRAVIS quite often: the subjects estimate the percentage of time they spent in EXTRAVIS at 70% on average. In itself, this percentage is meaningless: for example, in a related study it was observed that *"heavy use of a feature does not necessarily mean it (or the tool) helps to solve a task"*, and that *"repeated use may actually be a sign of frustration on the part of the user"* [40]. However, the questionnaire also shows that EXTRAVIS was used on seven of the eight tasks on average and that the tool was actually found useful in six of those tasks (86%). This is a strong indication that the Ecl+Ext subjects generally did not experience a resistance to using EXTRAVIS (resulting from, e.g., a poor understanding of the tool) *and* were quite successful in their attempts.

The latter assumption is further reinforced by the Ecl+Ext subjects' opinions on the speed and responsiveness of the tool, averaging a score of 1.35 on a scale of 0-2, which is between *"pretty OK: occasionally had to wait for information"* and *"very quickly: the information was shown"*

TABLE VI
DEBRIEFING QUESTIONNAIRE RESULTS

	<i>Eclipse</i>		<i>Ecl+Ext</i>	
	mean	stdev.	mean	stdev.
Miscellaneous				
Perceived time pressure (0-4)	2.18	1.19	2.06	0.66
Knowledge of dynamic analysis (0-4)	2.26	1.22	2.53	1.12
Perceived task difficulty (0-4)				
T1	1.00	0.71	1.65	0.79
T2.1	2.59	1.18	1.18	0.64
T2.2	2.24	1.15	1.53	0.80
T3.1	2.12	0.78	2.12	0.70
T3.2	2.29	0.92	1.53	0.72
T3.3	2.18	0.95	1.47	0.94
T4.1	2.40	0.63	2.65	0.86
T4.2	1.53	0.92	1.63	1.02
Average	2.04		1.72	
Use of EXTRAVIS				
No. of features used			7.12	2.67
No. of tasks conducted w/ tool			7.00	1.06
No. of tasks successfully conducted w/ tool			6.00	1.55
Est. % of time spent in tool			70.00	24.99
Perceived tool speed (0-2)			1.35	0.49

instantly". Furthermore, all 34 subjects turned out to be quite familiar with dynamic analysis: in the questionnaire they indicated an average knowledge level of 2.3 on a scale of 0-4 on this topic, which is between "I'm familiar with it and can name one or two benefits" and "I know it quite well and performed it once or twice".

As a side note, in a related study [48], no correlation could be identified between the subjects' experience levels and their performance. While in our experiment the same holds for the Ecl+Ext group and for correctness in the Eclipse group, there *does* exist a negative correlation between expertise and the time effort in the latter group: a high average expertise yielded lower time requirements, and vice versa. This observation partly underlines the importance of an adequate selection procedure when recruiting subjects for software engineering experiments.

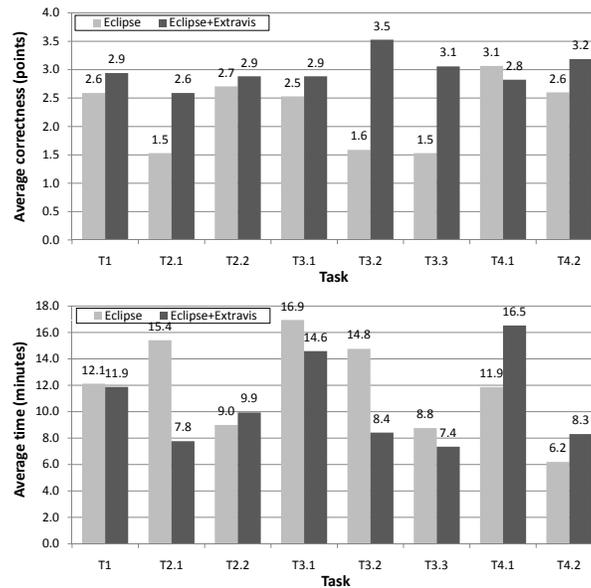


Fig. 4. Averages per task.

C. Individual task performance

To address our third research question, whether are certain types of comprehension tasks that benefit most from the use of EXTRAVIS (see Section III-A) we examine the performance per task in more detail. Figure 4 shows the average scores and time spent by each group from a task perspective.

While the experiment concerned only eight tasks, our data does suggest a negative correlation between time spent and correctness, in the sense that relatively little effort and a relatively high score (and vice versa) often go hand in hand.

Task T1

The goal of the first task was to identify and globally understand the most prominent stages in a typical CHECKSTYLE scenario (Table III). The groups scored equally well on this task and required similar amounts of time. According to the motivations of their solutions, the Eclipse group typically studied the `main()` method: however, such important phases as the building and parsing of an AST were often missing because they are not directly visible at the `main()` level. On the other hand, the EXTRAVIS users mostly studied an actual execution scenario through the

massive sequence view, which proved quite effective and led to slightly more accurate solutions.

Task T2.1

Task T2.1 concerned a fan-in/fan-out analysis that turned out to be significantly easier for the Ecl+Ext group, who scored 1.1 more points and needed only half the time. This is presumably explained by EXTRAVIS' circular view, from which *all* classes and their interrelationships can be directly interpreted. The Eclipse group mostly carried out a manual search for utility-like classes, opening numerous source files in the process, which is time-consuming and does not necessarily yield optimal results.

Task T2.2

This task was similar to the previous one, except that the focus was more on coupling. While there still exists a performance difference, it is much smaller this time round. According to the given solutions, the Ecl+Ext group again resorted to the circular view to look for high edge concentrations, while the Eclipse group mostly went searching for specific imports. The fact that a more specific (and automated) search was possible in this case may account for the improved performance of the latter group.

Task T3.1

Task T3.1 asked the participants to study a certain check to understand its life cycle, from creation to destruction. The performance difference here was quite subtle, with the Ecl+Ext group apparently having had a small advantage. Eclipse users typically studied the check's source code and started a more broad investigation from there. EXTRAVIS users mostly used our tool to highlight the check in the given execution trace and examine the interactions that were found. Interestingly, only a handful of subjects discovered that the checks are in fact dynamically loaded, and both groups often missed the explicit destruction of each check at the end of execution, which is not easily observed in Eclipse nor in EXTRAVIS.

Task T3.2

The goal of this follow-up task was to understand the protocol between a check and a certain key class, and asked the subjects to provide a list of interactions between these classes. The fact that the check at hand is an extension of a superclass that is an extension in itself, forced the Eclipse group to distribute its focus across each and every class in the check's type hierarchy.

EXTRAVIS users often highlighted the mutual interactions of the two classes at hand in the tool. As suggested by Figure 4, the latter approach is both faster and much more accurate (as there is a smaller chance of calls being missed).

Task T3.3

This task was similar to the previous one, except that it revolved around another type of check. The difference is that this check is dependent on the AST of the input source file, whereas the check in task T3.2 operates directly on the file. Finding the additional interactions was not too difficult for the EXTRAVIS users, who could follow a similar routine to last time. On the other hand, in Eclipse the subtle differences were often overlooked, especially if it was not understood that (and why) this check is fundamentally different from the previous one.

Task T4.1

Task T4.1 posed the challenging question of how CHECKSTYLE's error handling mechanism is implemented. It is the only task on which the Ecl+Ext group was clearly outperformed in terms of both time and correctness. The Eclipse group rated the difficulty of this task at 2.4, which is between "intermediate" and "difficult", whereas EXTRAVIS users rated the difficulty of this task at 2.65, leaning toward "difficult". An important reason might be that EXTRAVIS users did not know exactly *what to look for* in the execution trace, because the question was rather abstract in the sense that no clear starting point was given. On the other hand, the Eclipse group mostly used one of the checks as a baseline and followed the error propagation process from there. The latter approach is typically faster: the availability of EXTRAVIS may have been a distraction rather than an added value in this case.

Task T4.2

The focus in the final task was on testing the behavior of a check: given that a new check has been written and an input source file is available, how can we test if it works correctly? The Ecl+Ext group often searched the execution traces for communication between the check and the violation container class, which is quite effective once that class has been found. The Eclipse group had several choices. A few subjects tried to understand the check and apply this knowledge on the given input source file, i.e., understand which items the check is looking for, and then verify if these items occur in the input source file. Others tried to relate the check's

typical warning message (once it was determined) to example outputs given in the handouts; yet others used the Eclipse debugger, e.g., by inserting breakpoints or print statements in the error handling mechanism. With the exception of debugging, most of the latter approaches are quite time-consuming, if successful at all. Still, we observe no large difference in time spent: the fact that eight members of the Eclipse group had already exceeded the time limit at this point may have caused them to hurry, thereby reducing not only the time effort but also the score.

Summary

Following our interpretation of the individual task performance, we now formulate an analytical generalization [54] based on the quantitative results discussed earlier, the debriefing questionnaire results, and the four case studies from our earlier work [7].

Global structural insight. From the results of tasks T2.1 and T2.2 it has become clear that EXTRAVIS' circular view is of great help in grasping the structural relationships of the subject system. In particular, the bundling feature ensures that the many relations can all be shown simultaneously on a single screen. This poses a great advantage to using a standard IDE, in which often involves browsing through multiple files when a high-level structural insight is required. While any trace visualization technique could be helpful for such tasks, it should provide some means of visualizing the system's structural decomposition (e.g., UML sequence diagrams with hierarchically ordered lifelines [55]).

Global behavioral insight. In addition to structural insight, EXTRAVIS provides a navigable overview of an entire execution trace through the massive sequence view. As illustrated in earlier case studies and in task T1, this view visualizes the trace such that patterns can be visually distinguished. These patterns correspond to execution phases, the identification of which can be quite helpful in decomposing the subject system's behavior into smaller, more tractable pieces of functionality. In the case of CHECKSTYLE, this approach turned out to reveal more accurate information than could be derived from examining the `main()` method. A trace visualization technique must include some sort of navigable overview for it to be useful for such tasks.

Detailed behavioral insight. One of the main benefits of dynamic analysis is that occurrences of late binding are resolved, i.e., the maintainer can observe the actual objects involved in an execution scenario. This contributes to a more detailed understanding of a program's behavior. This is illustrated by tasks T3.2 and T3.3, which proved quite difficult for the Eclipse group

as these tasks concerned the identification of inherited methods, which are difficult to track down unless some form of run-time analysis is possible. We expect this particular advantage of dynamic analysis to be exploitable by any trace visualization technique.

Goal-oriented strategy. Trace visualization is not always the best solution: the results for task T4.1 showed a clear advantage for the Eclipse group. We believe that the reason can be generalized as follows: dynamic analysis typically involves a goal-oriented strategy, in the sense that one must know what to look for. (This follows from the fact that an appropriate execution scenario must be chosen.) If such a strategy is not feasible, e.g., because there is no clear starting point (such as the name of a certain class), then a strong reliance on dynamic analysis will result in mere confusion, which means that one must resort to traditional solutions such as the IDE instead.

D. Related experiments

There exist no earlier studies in the literature that offer quantitative evidence of the added value of trace visualization techniques for program comprehension. We therefore describe the experiments that are most closely related to our topic.

The aforementioned article from Bennett et al. concerned a user study involving a sequence diagram reconstruction tool [40]. Rather than measure its added value for program comprehension, they sought to characterize the *manner* in which the tool is used in practice. To this end, they had six subjects perform a series of comprehension tasks, and measured when and how the tool features were used. Among their findings was that tool features are not often formally evaluated in literature, and that heavily used tool features may indicate confusion among the users. Another important observation was that much time was spent on *scrolling*, which supports our hypothesis that EXTRAVIS saves time as it shows all information on a single screen.

Quante performed a controlled experiment to assess the benefits of Dynamic Object Process Graphs (DOPGs) for program comprehension [48]. While these graphs are built from run-time data, they do not actually visualize execution traces. The experiment involved 25 students and a series of feature location tasks for two subject systems. The use of DOPGs by his experimental group led to a significant decrease in time and a significant increase in correctness in case of the first system; however, the differences in case of the second system were *not* statistically significant. This suggests that evaluations on additional systems are also desirable for EXTRAVIS

and should be considered as future work. Also of interest is that the latter subject system was four times smaller than the former, but had three DOPGs associated with it instead of one. This may have resulted in an information overload on the part of the user, once more suggesting that users are best served by as little information as possible.

Among the contributions by Hamou-Lhadj and Lethbridge are encouraging quantitative results with respect to their trace summarization algorithm, effectively reducing large traces to a little as 0.5% of the original size [4]. However, the measurements performed relate to the effectiveness of the algorithm in terms of *reduction power*, rather than its added value in actual comprehension tasks.

VI. THREATS TO VALIDITY

This section discusses the validity threats in our experiment and the manners in which we have addressed them. We have identified three types of validity threats: (1) internal validity, referring to the cause-effect inferences made during the analysis; (2) external validity, concerning the generalizability of the results to different contexts; and (3) construct validity, seeking agreement between a theoretical concept and a specific measuring procedure.

A. Internal validity

Subjects. There exist several internal validity threats that relate to the subjects used in this experiment. First of all, the subjects may not have been sufficiently competent. We have reduced this threat through the a priori assessment of the subjects' competence in five relevant fields, which pointed out that all subjects had at least an elementary knowledge of Eclipse (2.47 in Figure 2) and no expert knowledge of CHECKSTYLE. Furthermore, participants could ask questions on both tools during the experiments, and a quick reference chart was available.

Second, their knowledge may not have been fairly distributed across the control group and experimental group. This threat was alleviated by grouping the subjects such that their expertise was evenly distributed across the groups (Figure 2).

Third, the subjects may not have been properly motivated, or may have had too much knowledge of the experimental goal. The former threat is mitigated by the fact that they all participated on a voluntary basis; as for the latter, the subjects were not familiar with the actual research questions or hypotheses (although they may have guessed).

Tasks. The comprehension tasks were designed by the authors of this paper, and therefore may have been biased toward EXTRAVIS (as this tool was also designed by the authors). To avoid this threat, we have applied an established task framework [51] to ensure that many aspects of typical comprehension contexts are covered. As a result, the tasks concerned both global and detailed knowledge, and both static and dynamic aspects.

Another task-related threat is that the tasks may have been too difficult. We refute this possibility on the basis of the correctness results, which show that maximum scores were occasionally awarded in both groups for all but one task (T3.1), which in the Eclipse group often yielded 3 points but never 4. However, the average scores for this task were a decent 2.53 (stdev. 0.51) and 2.88 (stdev. 0.86) in the Eclipse group and Ecl+Ext group, respectively. This point of view is further reinforced by the subjects' opinions on the task difficulties: the task they found hardest (T4.1) yielded good scores, being 3.07 (stdev. 1.10) for the Eclipse group and 2.82 (stdev. 0.81) for the Eclipse+Extravis group.

Also related to the tasks is the possibility that the subjects' solutions were graded incorrectly. This threat was reduced in our experiment by creating concept solutions in advance and by having CHECKSTYLE's lead developer review and refine them. This resulted in a solution model that clearly states the required elements (and corresponding points) for each task. Furthermore, to verify the soundness of the reviewing process, the first two authors of this paper independently reviewed the solutions of five random subjects: on each of the five occasions the difference was no higher than one point (out of the maximum of 32 points), suggesting a high inter-rater reliability.

Miscellaneous. The results may have been influenced by time constraints that were too loose or too strict. We have attempted to circumvent this threat by performing two pilot studies, which led to the removal of two tasks. Still, not all subjects finished the tasks in time, but the average time pressure (as indicated by the subjects in the debriefing questionnaire) was found to be 2.18 (stdev. 1.19) in the Eclipse group and 2.06 (stdev. 0.66) in the Ecl+Ext group on a scale of 0-4, which roughly corresponds to only a "*fair amount of time pressure*". Also, in our results analysis we have disregarded the last two tasks, upon which only three out of the 34 subjects still exceeded the time limit.

As several test subjects did not finish tasks T4.1 and T4.2 (within time), we decided to elimi-

nate these tasks from the analysis of our results. This removal may have benefited the EXTRAVIS results because task T4.1 is one of the few tasks at which the Eclipse group outperformed the EXTRAVIS users. Fortunately, with EXTRAVIS shown to be 43% more accurate and 21% less time-consuming, the conclusion that EXTRAVIS constitutes a significant added value for program comprehension would likely still be valid if tasks T4.1 and T4.2 were taken into account. Future refinements of the experimental design should examine optimizations of the time limit policy.

The two execution traces that we provided to the experimental group for use in EXTRAVIS are relatively small, containing 31,260 and 17,126 calls respectively. The fact that these traces are relatively small might influence the usability of EXTRAVIS: in particular, large traces could render EXTRAVIS a little less responsive and therefore a bit more time-consuming to use. However, earlier case studies [7] that we performed with EXTRAVIS (involving much larger traces) lead us to believe that the usability impact of using larger traces is probably minor.

Furthermore, our statistical analysis may not be completely accurate due to the missing data points that we mentioned in Section IV. This concerned two subjects who did not finish the last two tasks and one subject who did not finish the last task. Fortunately, the effect of the missing timing and correctness data points on our calculations is negligible: had the subjects finished the tasks, their total time spent and average score could have been higher, but this would only have affected the analysis of all eight tasks whereas our focus has been on the first six.

Another validity threat could be the fact that the control group only had access to the Eclipse IDE, whereas the experimental group also received two execution traces (next to Eclipse and the EXTRAVIS tool). However, we believe that the Eclipse group would not have benefited from the availability of execution traces because they are too large to be navigated without any tool support.

Lastly, it could be suggested that Eclipse is more powerful if additional plugins are used. However, as evidenced by the results of the debriefing questionnaire, only two subjects named specific plugins that would have made the tasks easier, and these related to only two of the eight tasks. We therefore expect that additional plugins would not have had a significant impact.

B. External validity

The generalizability of our results could be hampered by the limited representativeness of the subjects, the tasks, and CHECKSTYLE as a subject system.

Concerning the subjects, the use of professional developers instead of (mainly) Ph.D. candidates and M.Sc. students could have yielded different results. Unfortunately, motivating people from industry to sacrifice two hours of their precious time is quite difficult. Nevertheless, against the background of related studies that often employ undergraduate students, we assume the expertise levels of our 34 subjects to be relatively high. This assumption is partly reinforced by the (subjective) a priori assessment, in which the subjects rated themselves as being “*advanced*” with Java (avg. 3.06, stdev. 0.65), and “*regular*” at using Eclipse (avg. 2.47, stdev. 0.90). We acknowledge that our subjects’ knowledge of dynamic analysis may have been greater than in industry, averaging 2.26 (Table VI).

Another external validity threat concerns the comprehension tasks, which may not reflect real maintenance situations. We tried to neutralize this threat by relying on Pacione’s framework [51], which is based on activities often found in software visualization and the comprehension evaluation literature. The resulting tasks were reasonably complicated: Both groups encountered a task of which they rated the difficulty between 2.5 and 3.0, roughly corresponding to “difficult” (See the debriefing questionnaire results in Table VI). Furthermore, they also included an element of “surprise”: Task 3.1, for example, required the subjects to describe the life cycle of a given object, which made the majority of subjects enter in a fruitless search for its constructor, whereas the object was in fact dynamically loaded. Last but not least, the tasks concerned open questions, which approximate real life contexts better than multiple choice questions do. Nevertheless, arriving at a representative set of tasks that is suitable for use in experiments by different researchers is a significant challenge, which warrants further research.

Finally, the use of a different subject system (or additional runs) may have yielded different or more reliable results. CHECKSTYLE was chosen on the basis of several important criteria: in particular, finding another system of which the experimenters have sufficient knowledge is not trivial. Moreover, an additional case (or additional run) imposes twice the burden on the subjects or requires more of them. While this may be feasible in case the groups consist exclusively of students, it is not realistic in case of Ph.D. candidates or professional developers because they often have little time to spare.

C. Construct validity

In our experiment, we assessed the added value of our EXTRAVIS tool for program comprehension, and sought to generalize this added value to trace visualization techniques in general (Section V-C). However, it should be noted that the experiment does not enable a *distinction* between EXTRAVIS and trace visualization: we cannot tell whether the performance improvement should be attributed to trace visualization in general or to specific aspects of EXTRAVIS (e.g., the circular bundle view). To characterize the difference, there is a need for similar experiments involving other trace visualization techniques.

As another potential threat to construct validity, the control group did not have access to the execution traces. This may have biased the experimental group because they had more data to work with. The rationale behind this decision was our intent to mimic real-life working conditions, in which software engineers often limit themselves to the use of the IDE. The subjects could still study the behavior of the application using, e.g., the built-in debugger in Eclipse (which in the experiment was available to both groups and was indeed used by some).

VII. CONCLUSIONS

In this paper, we have reported on a controlled experiment that was aimed at the quantitative evaluation of EXTRAVIS, our tool for execution trace visualization. We designed eight typical tasks aimed at gaining an understanding of an open source program, and measured the performance of a control group (using the Eclipse IDE) and an experimental group (using both Eclipse and EXTRAVIS) in terms of time spent and correctness.

The results clearly illustrate EXTRAVIS' usefulness for program comprehension. With respect to time, the added value of EXTRAVIS was found to be statistically significant: on average, the EXTRAVIS group spent 22% less time on the given tasks. In terms of correctness, the results turned out even more convincing: EXTRAVIS' added value was again statistically significant, with the EXTRAVIS users scoring 43% more points on average. For the tasks that we considered, these results testify to EXTRAVIS' benefits compared to conventional tools: in this case, the Eclipse IDE.

To determine which types of tasks are best suited for EXTRAVIS or for trace visualization in general, we studied the group performance per task in more detail. While inferences drawn from one experiment and eight tasks cannot be conclusive, the experimental results do provide a

strong indication as to EXTRAVIS' strengths. First, questions that require insight into a system's structural relations are solved relatively easily due to EXTRAVIS' circular view, as it shows *all* of the system's structural entities and their call relationships on one single screen. Second, tasks that require a user to globally understand a system's behavior are easier to tackle when a visual representation of a trace is provided, as it decomposes the system's execution into tractable parts. Third, questions involving a detailed behavioral understanding seem to benefit greatly from the fact that dynamic analysis reveals the actual objects involved in each interaction, saving the user the effort of browsing multiple source files.

This paper demonstrates the potential of trace visualization for program comprehension, and paves the way for other researchers to conduct similar experiments. The work described in this paper makes the following contributions:

- A systematic literature survey of existing trace visualization techniques in the literature, and a description of the 21 contributions that were found.
- The design of a controlled experiment for the quantitative evaluation of trace visualization techniques for program comprehension, involving eight reusable tasks and a validated solution model.
- The execution of this experiment on a group of 34 representative subjects, demonstrating a 22% decrease in time effort and a 43% increase in correctness.
- A discussion on the types of tasks for which EXTRAVIS, and trace visualization in general, are best suited.

A. Future work

As mentioned in Section V-D, a related study has pointed out that results may differ quite significantly across different subject systems. It is therefore part of our future directions to replicate our experiment on another subject system.

Furthermore, we seek collaborations with fellow researchers to evaluate other existing trace visualization techniques. By subjecting other such techniques to the same experimental procedure, we might be able to quantify their added values for program comprehension as well, and compare their performance to that of EXTRAVIS.

Finally, we believe that strong quantitative results such as the ones presented in this study could play a crucial role in making industry realize the potential of dynamic analysis in their

daily work. In particular, they might be interested to incorporate trace visualization tools in their development cycle, and be willing to collaborate in a longitudinal study for us to investigate the long-term benefits of dynamic analysis in practice. Another aim of such a longitudinal study could be to shed light on how software engineers using a dynamic analysis tool define an execution scenario, how often they do this, and how much time they spend on it.

ACKNOWLEDGMENTS

This research is sponsored by NWO via the Jacquard Reconstructor project. We would like to thank the 34 subjects for their participation, Danny Holten for his implementation of EXTRAVIS, Cathal Boogerd for performing one of the pilot studies, and Bart Van Rompaey for assisting in the experimental design. Also, many thanks to CHECKSTYLE's lead developer, Oliver Burn, who assisted in the design of our task review protocol.

REFERENCES

- [1] T. A. Corbi, "Program understanding: Challenge for the 1990s," *IBM Systems Journal*, vol. 28, no. 2, pp. 294–306, 1989.
- [2] V. R. Basili, "Evolving and packaging reading technologies," *Journal of Systems & Software*, vol. 38, no. 1, pp. 3–12, 1997.
- [3] S. P. Reiss and M. Renieris, "Encoding program executions," in *Proc. International Conference on Software Engineering (ICSE)*, pp. 221–230, IEEE C.S., 2001.
- [4] A. Hamou-Lhadj and T. C. Lethbridge, "Summarizing the content of large traces to facilitate the understanding of the behaviour of a software system," in *Proc. International Conference on Program Comprehension (ICPC)*, pp. 181–190, IEEE C.S., 2006.
- [5] W. De Pauw, R. Helm, D. Kimelman, and J. M. Vlissides, "Visualizing the behavior of object-oriented systems," in *Proc. Conference on Object-Oriented Programming Systems, Languages & Applications (OOPSLA)*, pp. 326–337, ACM, 1993.
- [6] D. F. Jerding, J. T. Stasko, and T. Ball, "Visualizing interactions in program executions," in *Proc. International Conference on Software Engineering (ICSE)*, pp. 360–370, ACM, 1997.
- [7] B. Cornelissen, A. Zaidman, D. Holten, L. Moonen, A. van Deursen, and J. J. van Wijk, "Execution trace analysis through massive sequence and circular bundle views," *Journal of Systems & Software*, vol. 81, no. 11, pp. 2252–2268, 2008.
- [8] B. Cornelissen, A. Zaidman, B. Van Rompaey, and A. van Deursen, "Trace visualization for program comprehension: A controlled experiment," in *Proc. Int. Conf. on Program Comprehension (ICPC)*, IEEE C.S., 2009. To appear.
- [9] B. Cornelissen, A. Zaidman, A. van Deursen, L. Moonen, and R. Koschke, "A systematic survey of program comprehension through dynamic analysis," *IEEE Transactions on Software Engineering*, vol. 35, no. 5, pp. 684–702, 2009.
- [10] A. Zaidman and S. Demeyer, "Managing trace data volume through a heuristical clustering process based on event execution frequency," in *Proc. European Conference on Software Maintenance & Reengineering (CSMR)*, pp. 329–338, IEEE C.S., 2004.

- [11] M. F. Kleyn and P. C. Gingrich, "Graphtrace - understanding object-oriented systems using concurrently animated views," in *Proc. Conference on Object-Oriented Programming Systems, Languages & Applications (OOPSLA)*, pp. 191–205, ACM, 1988.
- [12] W. De Pauw, D. Lorenz, J. Vlissides, and M. Wegman, "Execution patterns in object-oriented visualization," in *Proc. USENIX Conference on Object-Oriented Technologies & Systems (COOTS)*, pp. 219–234, USENIX, 1998.
- [13] W. De Pauw, E. Jensen, N. Mitchell, G. Sevitsky, J. M. Vlissides, and J. Yang, "Visualizing the execution of Java programs," in *Proc. ACM 2001 Symposium on Software Visualization (SOFTVIS)*, pp. 151–162, ACM, 2001.
- [14] W. De Pauw, S. Krasikov, and J. F. Morar, "Execution patterns for visualizing web services," in *Proc. ACM 2006 Symposium on Software Visualization (SOFTVIS)*, pp. 37–45, ACM, 2006.
- [15] K. Koskimies and H. Mössenböck, "Scene: Using scenario diagrams and active text for illustrating object-oriented programs," in *Proc. International Conference on Software Engineering (ICSE)*, pp. 366–375, IEEE C.S., 1996.
- [16] D. F. Jerding and S. Rugaber, "Using visualization for architectural localization and extraction," in *Proc. Working Conference on Reverse Engineering (WCRE)*, pp. 56–65, IEEE C.S., 1997.
- [17] T. Systä, "On the relationships between static and dynamic models in reverse engineering Java software," in *Proc. Working Conference on Reverse Engineering (WCRE)*, pp. 304–313, IEEE C.S., 1999.
- [18] T. Systä, K. Koskimies, and H. A. Müller, "Shimba: an environment for reverse engineering Java software systems," *Software: Practice & Experience*, vol. 31, no. 4, pp. 371–394, 2001.
- [19] T. S. Souder, S. Mancoridis, and M. Salah, "Form: A framework for creating views of program executions," in *Proc. International Conference on Software Maintenance (ICSM)*, pp. 612–620, IEEE C.S., 2001.
- [20] R. Oechsle and T. Schmitt, "JAVAVIS: Automatic program visualization with object and sequence diagrams using the Java Debug Interface (JDI)," in *Proc. ACM 2001 Symposium on Software Visualization (SOFTVIS)*, pp. 176–190, ACM, 2001.
- [21] A. Hamou-Lhadj and T. C. Lethbridge, "Compression techniques to simplify the analysis of large execution traces," in *Proc. International Workshop on Program Comprehension (IWPC)*, pp. 159–168, IEEE C.S., 2002.
- [22] A. Hamou-Lhadj, T. C. Lethbridge, and L. Fu, "Challenges and requirements for an effective trace exploration tool," in *Proc. International Workshop on Program Comprehension (IWPC)*, pp. 70–78, IEEE C.S., 2004.
- [23] A. Hamou-Lhadj, E. Braun, D. Amyot, and T. C. Lethbridge, "Recovering behavioral design models from execution traces," in *Proc. European Conference on Software Maintenance and Reengineering (CSMR)*, pp. 112–121, IEEE C.S., 2005.
- [24] M. Salah and S. Mancoridis, "Toward an environment for comprehending distributed systems," in *Proc. Working Conference on Reverse Engineering (WCRE)*, pp. 238–247, IEEE C.S., 2003.
- [25] M. Salah and S. Mancoridis, "A hierarchy of dynamic software views: From object-interactions to feature-interactions," in *Proc. International Conference on Software Maintenance (ICSM)*, pp. 72–81, IEEE C.S., 2004.
- [26] M. Salah, T. Denton, S. Mancoridis, A. Shokoufandeh, and F. I. Vokolos, "Scenariographer: A tool for reverse engineering class usage scenarios from method invocation sequences," in *Proc. International Conference on Software Maintenance (ICSM)*, pp. 155–164, IEEE C.S., 2005.
- [27] M. Salah, S. Mancoridis, G. Antoniol, and M. Di Penta, "Scenario-driven dynamic analysis for comprehending large software systems," in *Proc. European Conference on Software Maintenance & Reengineering (CSMR)*, pp. 71–80, IEEE C.S., 2006.
- [28] L. C. Briand, Y. Labiche, and Y. Miao, "Towards the reverse engineering of UML sequence diagrams," in *Proc. Working Conference on Reverse Engineering (WCRE)*, pp. 57–66, IEEE C.S., 2003.

- [29] L. C. Briand, Y. Labiche, and J. Leduc, "Tracing distributed systems executions using AspectJ," in *Proc. International Conference on Software Maintenance (ICSM)*, pp. 81–90, IEEE C.S., 2005.
- [30] L. C. Briand, Y. Labiche, and J. Leduc, "Toward the reverse engineering of UML sequence diagrams for distributed Java software," *IEEE Transactions on Software Engineering*, vol. 32, no. 9, pp. 642–663, 2006.
- [31] A. Kuhn and O. Greevy, "Exploiting the analogy between traces and signal processing," in *Proc. International Conference on Software Maintenance (ICSM)*, pp. 320–329, IEEE C.S., 2006.
- [32] O. Greevy, M. Lanza, and C. Wyseier, "Visualizing live software systems in 3D," in *Proc. ACM 2006 Symposium on Software Visualization (SOFTVIS)*, pp. 47–56, ACM, 2006.
- [33] J. Koskinen, M. Kettunen, and T. Systä, "Profile-based approach to support comprehension of software behavior," in *Proc. International Conference on Program Comprehension (ICPC)*, pp. 212–224, IEEE C.S., 2006.
- [34] S. Simmons, D. Edwards, N. Wilde, J. Homan, and M. Groble, "Industrial tools for the feature location problem: an exploratory study," *Journal of Software Maintenance & Evolution*, vol. 18, no. 6, pp. 457–474, 2006.
- [35] B. Cornelissen, A. van Deursen, L. Moonen, and A. Zaidman, "Visualizing testsuites to aid in software understanding," in *Proc. European Conference on Software Maintenance & Reengineering (CSMR)*, pp. 213–222, IEEE C.S., 2007.
- [36] R. Voets, "JRET: A tool for the reconstruction of sequence diagrams from program executions," Master's thesis, Delft University of Technology, 2008.
- [37] S. P. Reiss, "Visual representations of executing programs," *Journal of Visual Languages & Computing*, vol. 18, no. 2, pp. 126–148, 2007.
- [38] J. Jiang, J. Koskinen, A. Ruokonen, and T. Systä, "Constructing usage scenarios for API redocumentation," in *Proc. International Conference on Program Comprehension (ICPC)*, pp. 259–264, IEEE C.S., 2007.
- [39] B. Cornelissen, D. Holten, A. Zaidman, L. Moonen, J. J. van Wijk, and A. van Deursen, "Understanding execution traces using massive sequence and circular bundle views," in *Proc. International Conference on Program Comprehension (ICPC)*, pp. 49–58, IEEE C.S., 2007.
- [40] C. Bennett, D. Myers, D. Ouellet, M.-A. Storey, M. Salois, D. German, and P. Charland, "A survey and evaluation of tool features for understanding reverse engineered sequence diagrams," *Journal of Software Maintenance & Evolution*, vol. 20, no. 4, pp. 291–315, 2008.
- [41] A. R. Dalton and J. O. Hallstrom, "A toolkit for visualizing the runtime behavior of TinyOS applications," in *Proc. International Conference on Program Comprehension (ICPC)*, pp. 43–52, IEEE C.S., 2008.
- [42] M. J. Pacione, M. Roper, and M. Wood, "Comparative evaluation of dynamic visualisation tools," in *Proc. Working Conference on Reverse Engineering (WCRE)*, pp. 80–89, IEEE C.S., 2003.
- [43] A. Hamou-Lhadj and T. C. Lethbridge, "A survey of trace exploration tools and techniques," in *Proc. Conference of the Centre for Advanced Studies on Collaborative Research (CASCON)*, pp. 42–55, IBM Press, 2004.
- [44] N. Wilde and M. C. Scully, "Software Reconnaissance: Mapping program features to code," *Journal of Software Maintenance*, vol. 7, no. 1, pp. 49–62, 1995.
- [45] D. F. Jerding and J. T. Stasko, "The information mural: A technique for displaying and navigating large information spaces," *IEEE Transactions on Visualization & Computer Graphics*, vol. 4, no. 3, pp. 257–271, 1998.
- [46] M.-A. Storey, "Theories, methods and tools in program comprehension: past, present and future," in *Proc. International Workshop on Program Comprehension (IWPC)*, pp. 181–191, IEEE C.S., 2005.
- [47] C. F. J. Lange and M. R. V. Chaudron, "Interactive views to improve the comprehension of UML models - an experimental validation," in *Proc. International Conference on Program Comprehension (ICPC)*, pp. 221–230, IEEE C.S., 2007.

- [48] J. Quante, "Do dynamic object process graphs support program understanding? – a controlled experiment," in *Proc. International Conference on Program Comprehension (ICPC)*, pp. 73–82, IEEE C.S., 2008.
- [49] A. Zaidman, B. Van Rompaey, S. Demeyer, and A. van Deursen, "Mining software repositories to study co-evolution of production & test code," in *Proc. International Conference on Software Testing (ICST)*, pp. 220–229, IEEE C.S., 2008.
- [50] B. Van Rompaey and S. Demeyer, "Estimation of test code changes using historical release data," in *Proc. Working Conference on Reverse Engineering (WCRE)*, pp. 269–278, IEEE C.S., 2008.
- [51] M. J. Pacione, M. Roper, and M. Wood, "A novel software visualisation model to support software comprehension," in *Proc. Working Conference on Reverse Engineering (WCRE)*, pp. 70–79, IEEE C.S., 2004.
- [52] M. Di Penta, R. E. K. Stirewalt, and E. Kraemer, "Designing your next empirical study on program comprehension," in *Proc. International Conference on Program Comprehension (ICPC)*, pp. 281–285, IEEE C.S., 2007.
- [53] C. Wohlin, P. Runeson, M. Höst, M. C. Ohlsson, B. Regnell, and A. Wesslen, *Experimentation in software engineering - an introduction*. Kluwer Acad. Publ., 2000.
- [54] R. K. Yin, *Case Study Research: Design and Methods*. Sage Publications Inc., 2003.
- [55] C. Riva and J. V. Rodriguez, "Combining static and dynamic views for architecture reconstruction," in *Proc. European Conference on Software Maintenance & Reengineering (CSMR)*, pp. 47–55, IEEE C.S., 2002.

Subjects & Results

TABLE VII

CHARACTERIZATION OF THE SUBJECTS AND THEIR PERFORMANCE, ORDERED BY AVERAGE EXPERTISE.

Subj.	Affil.	M.Sc.	Avg	Expertise (0-4)					Performance (T1-T3)	
				Java	Ecl.	Chkst.	Lang.tech.	Rev.eng.	Time	Correctness
1	SIG	VU	3.4	4	4	1	4	4	54	22
2	CWI	UvA	3.2	4	4	1	4	3	56	7
3	TUD	.at	3.0	4	4	0	3	4	67	13
4	SIG	UU	3.0	3	3	3	3	3	55	14
5	TUD	TUD	2.8	4	3	2	2	3	70	18
6	SIG	TU/e	2.8	4	4	1	2	3	68	17
7	TUD	UU	2.8	4	4	0	4	2	68	20
8	CWI	UvA	2.8	4	3	3	2	2	38	16
9	SIG	.pt	2.4	3	3	1	2	3	64	18
10	SIG	.pt	2.4	3	3	1	2	3	84	9
11	TUD	KUN	2.2	3	3	1	3	1	59	18
12	TUD	.pt	2.2	3	2	0	3	3	37	11
13	TUD	TUD	2.2	3	2	0	3	3	63	15
14	.pt	.pt	2.2	3	3	2	2	1	78	8
15	.fi	.fi	2.0	4	1	1	0	4	70	22
16	UT	UT	2.0	3	1	1	3	2	70	14
17	TUD	TUD	2.0	2	3	0	3	2	96	9
18	TUD	TUD	2.0	3	3	0	2	2	70	13
19	TUD	.fr	2.0	3	2	0	2	3	81	19
20	.be	.be	2.0	3	2	0	2	3	84	14
21	.au	.au	2.0	2	2	0	4	2	77	15
22	TUD	TUD	1.8	3	2	0	2	2	66	19
23	TUD	TUD	1.8	3	2	0	3	1	36	19
24	TUD	.es	1.8	3	2	0	2	2	58	19
25	TUD	RUG	1.8	3	2	0	3	1	38	20
26	TUD	TUD	1.8	3	3	0	2	2	72	22
27	TUD	TUD	1.8	3	2	0	3	0	97	14
28	TUD	TUD	1.6	3	2	0	2	1	70	18
29	TUD	TUD	1.6	2	1	0	2	2	44	22
30	TUD	UU	1.6	2	1	0	3	2	88	15
31	TUD	TUD	1.4	3	2	0	1	1	69	13
32	TUD	.de	1.4	3	2	0	2	0	79	7
33	.be	.be	1.2	2	2	0	1	1	102	5
34	TUD	TUD	1.0	2	1	0	2	0	100	11

Handouts

Introduction

Thank you for your willingness to participate in this experiment! Empirical studies are not very common in the field of software understanding because this field has a strong cognitive aspect that is difficult to measure. This makes controlled experiments (such as the one in which you are now participating) all the more valuable. I hope you will find it an interesting experience.

The context of this experiment concerns a (fictive) developer who is asked to perform certain maintenance tasks on a system, but who is unfamiliar with its implementation. The focus of the experiment is not on *performing* these maintenance tasks, but rather on gaining the necessary *knowledge* and measuring the effort that is involved therein.

The case study in this experiment is Checkstyle. From the Checkstyle site:

Checkstyle is a development tool to help programmers write Java code that adheres to a coding standard. It automates the process of checking Java code to spare humans of this boring (but important) task. This makes it ideal for projects that want to enforce a coding standard.

In short, Checkstyle takes as inputs a Java source file, and an XML configuration file that specifies the coding standards that must be enforced, i.e., the *checks* that are to be used.

Most people are not familiar with Checkstyle's implementation. However, IDEs (such as Eclipse) and effective tools may be able to assist in understanding Checkstyle's inner workings, and most of the source code is fairly documented.

You are given 90 minutes for four comprehension tasks, which have been structured according to their maintenance contexts. Each task involves several related subtasks, some of which designed to help you on your way. The task ordering implies a top-down approach: we start by building a general knowledge of the system and then drill down to a more detailed understanding. For each of the subtasks, you are asked to write down the following items:

- Your answer.
- A motivation of your answer. In most cases it suffices to briefly describe how your answer was obtained.
- The *time* at which you started this subtask (important!).

Furthermore, you are asked (1) to not consult any other participants, (2) to perform the tasks in the order specified, and (3) to *not return to earlier tasks* because it affects the timing. Finally, while there is an online documentation available for Checkstyle, you are kindly requested *not* to use this, because we want to simulate a real life case in which up-to-date documentation is often lacking. Using the Internet is allowed only for Java-related resources (e.g., APIs).

We start off by describing the tools at your disposal. You are then presented with the comprehension tasks, at which point your 90 minutes start ticking. The experiment is concluded with a short questionnaire.

Using the Eclipse IDE

Eclipse is the IDE that you will be using to perform the tasks. You are expected to have a basic knowledge of Eclipse, but a quick “reference chart” is provided nonetheless. While this chart only shows the basic features, of course you are encouraged to use more advanced functionalities if you are familiar with them.

Your Eclipse setup contains a project with Checkstyle’s source (and links to its external libraries). Several aspects are worth mentioning:

- Checkstyle’s testsuite is not available to you. This reflects real life cases in which the testsuite is not complete, out of date, or non-existent at all.
- The experiment will not be concentrating on Checkstyle’s GUI.
- You may compile and run the application if so desired.

Finally, you have at your disposal an input source file, `Simple.java`; and two different configuration XML-files, `many_checks.xml` and `several_checks.xml`. The resulting outputs of running Checkstyle with these inputs are given on the next few pages.

Should you have trouble using Eclipse, please refer to the reference chart, or consult me (Bas).

Using Extravis

In addition to the Eclipse IDE, you will also have access to Extravis during the experiment. Extravis is a dynamic analysis tool, which means it provides information on the system's runtime behavior. In this experiment, this is done through *execution traces*, which were obtained by instrumenting Checkstyle and then having it run a certain execution scenario. Such traces contain:

- 1) A chronological ordering of all *method* and *constructor calls* that occurred during execution. Typically this amounts to thousands or even millions of events for each scenario.
- 2) The *actual* class instances (objects) on which these methods and constructors were invoked. This means that (e.g.) if class A inherits method `a()` from some superclass B, the trace will show the receiver of the call `A.a()` to be A, not B.

Extravis visualizes these execution traces and the program's package decomposition, and provides means to navigate this information.

A quick reference chart of Extravis has been provided as part of the handouts. In addition to depicting all method and constructor calls that occurred in the scenario, Extravis also shows the *actual parameters* and *actual return values* for those calls. The developer is thus provided with a rich and accurate source of information with respect to the scenario at hand.

You have two execution traces at your disposal:

`simple-many_checks` and `simple-several_checks`. The former trace is the result of Checkstyle's execution with `Simple.java` and `many_checks.xml` as its parameters; the latter trace was also obtained using `Simple.java`, but with `several_checks.xml` as the configuration file.

Finally, some aspects that are worth noting:

- Only Checkstyle's *core* functionalities were instrumented, which means that the resulting traces do not contain calls to, or from, external libraries or the JDK.
- Extravis provides two *linked* views: changes made in the one view are propagated toward the other.
- The leftmost view concentrates on visualizing the target system's structure and its (runtime!) interrelationships, whereas the rightmost view focuses more on navigating the trace.

In other words, Extravis answers questions related to a program's actual execution, and aims to provide insight in the interactions that take place.

You are free to use Extravis during your tasks whenever you see fit. Should you have trouble using the tool, please refer to the reference chart, or consult me (Bas).

Checkstyle outputs

```
bas@laptop:~/checkstyle-4.4$ java
    -jar checkstyle-all-4.4.jar
    -c several_checks.xml Simple.java
Starting audit...
Simple.java:11: Line has trailing spaces.
Simple.java:17:1: Line contains a tab character.
Simple.java:18:1: Line contains a tab character.
Simple.java:19:1: Line contains a tab character.
Simple.java:20:1: Line contains a tab character.
Simple.java:21:1: Line contains a tab character.
Simple.java:23:1: Line contains a tab character.
Simple.java:24:1: Line contains a tab character.
Simple.java:25:1: Line contains a tab character.
Simple.java:26:1: Line contains a tab character.
Simple.java:27:1: Line contains a tab character.
Simple.java:28:1: Line contains a tab character.
Simple.java:29: Line has trailing spaces.
Simple.java:29:1: Line contains a tab character.
Simple.java:30:1: Line contains a tab character.
Simple.java:31:1: Line contains a tab character.
Simple.java:32:1: Line contains a tab character.
Simple.java:33:1: Line contains a tab character.
Simple.java:34:1: Line contains a tab character.
Simple.java:35:1: Line contains a tab character.
Simple.java:36:1: Line contains a tab character.
Simple.java:37:1: Line contains a tab character.
Simple.java:38:1: Line contains a tab character.
Simple.java:40:1: Line contains a tab character.
Simple.java:40:9: Missing a Javadoc comment.
Simple.java:41:1: Line contains a tab character.
Simple.java:42:1: Line contains a tab character.
Simple.java:43:1: Line contains a tab character.
Simple.java:44:1: Line contains a tab character.
```

```
Simple.java:45:1: Line contains a tab character.  
Simple.java:46:1: Line contains a tab character.  
Simple.java:47:1: Line contains a tab character.  
Audit done.
```

```
bas@laptop:~/checkstyle-4.4$ java  
    -jar checkstyle-all-4.4.jar  
    -c many_checks.xml Simple.java  
Starting audit...  
/home/bas/checkstyle-4.4/package.html:0: Missing  
    package documentation file.  
Simple.java:11: Line has trailing spaces.  
Simple.java:17:1: Line contains a tab character.  
Simple.java:17:9: Missing a Javadoc comment.  
Simple.java:18:1: Line contains a tab character.  
Simple.java:18:9: Missing a Javadoc comment.  
Simple.java:19:1: Line contains a tab character.  
Simple.java:19:9: Missing a Javadoc comment.  
Simple.java:19:23: '<' is not preceded with  
    whitespace.  
Simple.java:19:24: '<' is not followed by  
    whitespace.  
Simple.java:19:31: '>' is not preceded with  
    whitespace.  
Simple.java:20:1: Line contains a tab character.  
Simple.java:20:9: Missing a Javadoc comment.  
Simple.java:21:1: Line contains a tab character.  
Simple.java:21:9: Missing a Javadoc comment.  
Simple.java:23:1: Line contains a tab character.  
Simple.java:24:1: Line contains a tab character.  
Simple.java:25:1: Line contains a tab character.  
Simple.java:26:1: Line contains a tab character.  
Simple.java:27:1: Line contains a tab character.  
Simple.java:28:1: Line contains a tab character.
```

```
Simple.java:29: Line has trailing spaces.
Simple.java:29:1: Line contains a tab character.
Simple.java:30:1: Line contains a tab character.
Simple.java:31:1: Line contains a tab character.
Simple.java:32:1: Line contains a tab character.
Simple.java:33:1: Line contains a tab character.
Simple.java:34:1: Line contains a tab character.
Simple.java:34:21: Parameter map should be final.
Simple.java:35:1: Line contains a tab character.
Simple.java:36:1: Line contains a tab character.
Simple.java:37:1: Line contains a tab character.
Simple.java:38:1: Line contains a tab character.
Simple.java:40:1: Line contains a tab character.
Simple.java:40:9: Method 'initialize' is not
                 designed for extension - needs
                 to be abstract, final or empty.
Simple.java:40:9: Missing a Javadoc comment.
Simple.java:41:1: Line contains a tab character.
Simple.java:41:38: '<' is not preceded with
                 whitespace.
Simple.java:41:39: '<' is not followed by
                 whitespace.
Simple.java:41:46: '>' is not preceded with
                 whitespace.
Simple.java:41:47: '>' is not followed by
                 whitespace.
Simple.java:42:1: Line contains a tab character.
Simple.java:43:1: Line contains a tab character.
Simple.java:44:1: Line contains a tab character.
Simple.java:45:1: Line contains a tab character.
Simple.java:46:1: Line contains a tab character.
Simple.java:47:1: Line contains a tab character.
Audit done.
```

Tasks

1. Gaining a general understanding

The first thing the developer might desire is a first impression of how Checkstyle works, especially in case domain knowledge is lacking.

- **Task 1.** [current time: . . . : . . .]

Having glanced through the available information for several minutes, which do you think are the main stages in a typical (non-GUI) Checkstyle scenario? Formulate your answer from a *high-level* perspective: refrain from using identifier names and stick to a maximum of six main stages.

2. Identifying refactoring opportunities

In certain cases (not necessarily Checkstyle) it is desirable to modify the program's package hierarchy. Examples include the movement of tightly coupled classes to the same package, and the movement of classes with high fan-in and (almost) no fan-out to a utility package.

The *fan-in* of a class is defined as the number of distinct methods/constructors directed toward that class, not counting self-calls. Its *fan-out* is defined as the number of distinct methods/constructors directed toward other classes.

- **Task 2.1.** [current time: . . . : . . .]

Name three classes in Checkstyle that have a high fan-in and (almost) no fan-out.

Assume that a tight coupling is characterized by a relatively large number of *different* method calls between two structural entities (e.g., classes or packages).

- **Task 2.2.** [current time: . . . : . . .]

Name a class in the default package (i.e., classes not in any package) that could be a candidate for movement to the `api` package because of its tight coupling with classes therein.

3. Understanding the checking process

Checkstyle's purpose is the application of *checks* on its input source file. These checks each have their own class and are located in the `checks`-package. They can be written by a developer and contributed to the Checkstyle package: For example, one could write a check to impose a

limit on the number of methods in a class. If our developer wants to write a new check, one way to gain the necessary knowledge is to study existing checks (i.e., learning by example).

Let's assume that we want to know how `checks whitespace.TabCharacterCheck` interacts with the rest of the program, and that this check is part of the current configuration (and will therefore be applied).

- **Task 3.1.** [current time:]

In *general* terms, describe the life cycle of this check during execution: when is it created, what does it do and on whose command, and how does it end up?

Do not go into details yet, and use no more than five sentences.

The `TreeWalker` class plays an important role in Checkstyle's inner workings and interacts extensively with the checks. We now take a closer look at the protocol between `TreeWalker` and the various checks.

- **Task 3.2.** [current time:]

List the identifiers of all method/constructor calls that typically occur between `TreeWalker` and a `checks whitespace.TabCharacterCheck` instance, and the order in which they are called. Make sure you also take inherited methods/constructors into account.

- **Task 3.3.** [current time:]

In comparison to the calls listed in Task 3.2., which *additional* calls occur between `TreeWalker` and `checks coding.IllegalInstantiationCheck`? Can you think of a reason for the difference?

4. Understanding the violation reporting process

Once the developer has written a new check, he/she would like to know if it works, i.e., whether it reports warnings when appropriate. Consider the situation in which some check has encountered a violation.

- **Task 4.1.** [current time:]

How is the check's warning handled, i.e., where/how does it originate, how is it internally represented, and how is it ultimately communicated to the user?

Verifying whether a check actually found violations is not trivial: most of the warnings that are reported in Checkstyle's output (provided a few pages back) cannot be traced back directly to

the checks from which those warnings originate. Some reported warnings may even be quite confusing.

- **Task 4.2.** [current time: . . : . .]

Given `Simple.java` as the input source and `many_checks.xml` as the configuration, does `checks.whitespace.WhitespaceAfterCheck` report warnings? Specify how your answer was obtained.

Debriefing Questionnaire [current time: . . .]

The experiment is concluded with a short questionnaire in which we ask for your opinions on several experimental aspects. You may fill in your answers on the handouts themselves.

- On a scale of 1 to 5, how did you feel about the time pressure?
 - 1) too much time pressure: could not cope with it, regardless of task difficulty
 - 2) fair amount of time pressure: could certainly have done better with more time
 - 3) not so much time pressure: hurried a bit, but it was OK
 - 4) very little time pressure: felt quite comfortable
 - 5) no time pressure at all
- Regardless of the time given, how difficult would you rate the tasks? Please mark the appropriate difficulty for each of the tasks:

	impossible	difficult	intermediate	simple	trivial
Task 1					
Task 2.1					
Task 2.2					
Task 3.1					
Task 3.2					
Task 3.3					
Task 4.1					
Task 4.2					

- Which particular Eclipse features did you frequently use?
 - Package explorer
 - Open declaration
 - Open type hierarchy
 - Open call hierarchy
 - Text search
 - other:
 - ...
 - ...

- Do you feel that additional Eclipse plugins (that you know of) could have helped during the experiment? If so, please name those plugins and briefly explain how they would have assisted you.
...
...
- “Dynamic analysis” is the study of a program through its run-time behavior. Are you familiar with dynamic analysis and its benefits?
 - 1) Never heard of it
 - 2) I know what it is, more or less
 - 3) I’m familiar with it and could name one or two benefits
 - 4) I know it quite well and performed it once or twice
 - 5) I’ve used it on multiple occasions

And finally, several questions on your use of Extravis:

- During which tasks did you use Extravis, and in which of these tasks did you actually find it helpful?

	used it	used it successfully
Task 1		
Task 2.1		
Task 2.2		
Task 3.1		
Task 3.2		
Task 3.3		
Task 4.1		
Task 4.2		

- Can you give a rough estimate of the percentage of time that you spent on using Extravis?
...
- Which particular Extravis features did you use more than once? Please mark those features on your reference chart.
- Which (types of) Extravis features did you feel were missing?
 - Compare multiple traces side-by-side
 - Interactiveness in terms of the input (e.g., support for creating and visualizing own traces)
 - More readable / intuitive views of detailed interactions
 - Direct link from actual calls to their source code locations
 - Search capabilities
 - other:
...
...
- On a scale of 1 to 3, how did Extravis perform in terms of speed and responsiveness?
 - 1) Quite sluggish: I got impatient very often
 - 2) Pretty OK: occasionally had to wait for information

3) Very quickly: the information was shown instantly

- Did you experience a certain “resistance” to using Extravis and, instead, stuck to Eclipse as much as possible? If so, how would you explain this tendency? (multiple answers possible)
 - Time pressure
 - I’m so comfortable or skilled with Eclipse that I prefer to use Eclipse whenever I can
 - I felt that the tasks simply did not require more than Eclipse / source code
 - I’m not sufficiently familiar with Extravis (tutorial was too short)
 - I’m not (sufficiently) familiar with the benefits of using run-time information in general
 - Extravis is standalone rather than embedded in Eclipse, and I’m not comfortable with context switching
 - other:
 - ...
 - ...

Answer model

Task 1

The following stages largely capture a typical Checkstyle scenario, and represent the minimum for this question. We can be somewhat flexible since the task deals from a very high level perspective.

Assign one point for each stage that is contained by the given answer. Additional stages are permitted but do not yield additional points. Motivations are not necessary if the answer is meaningfully-named or self-explanatory.

- Initialization: command line parsing, or config reading, or environment setup (creation of checkers, listeners etc.)
- Source parsing: source input file is read/parsed, or AST construction
- Checking: input file is checked, or AST traversal
- Error reporting/Termination/results logging: conveyance of warnings, and teardown of application

Task 2.1

The second list below is alphabetically ordered and shows all classes of which the fanin is higher than the fanout. (Classes not in this list are obviously incorrect.) We were looking for classes with a (relatively) high fanin and a low fanout, so award points in case of appropriate proportions between the two. Examples include fanin 5 + fanout 0, fanin 10 + fanout 1, fanin 15 + fanout 2, etc. (flexible scale).

Each correct class receives one point; award four points in case all classes have appropriate fanins and fanouts of 0 or 1. In case no plausible motivation is given, award no points.

Task 2.2

The first list below shows all classes in the default package, and their degrees of coupling with the api package. The coupling values in this list were statically derived; it is defined as the sum of the no. of distinct calls from, and the no. of calls to, classes in the api package.

TreeWalker and Checker receive 4 points because in the context of the question (i.e., our definition of "coupling") they have the strongest api-coupling by far. The next four classes are awarded 2 points; all others receive none.

For certain alternative classes it can be argued that they have a strong coupling with the api because they communicate exclusively with the api (and with external libraries). While this is reasonable, it is not according to the given definition of coupling, and therefore receives only 2 points. Award one point in case one went looking in the api-package.

Task 3.1

The following elements largely capture the lifecycle of a TabCharacterCheck instance during the specified execution scenario. Again we can be a bit flexible due to the high-level perspective, but the elements below must be mentioned because they answer the four implicit sub-questions (where does it originate, what does it do, on whose command, how does it end up). Assign one point for each element that is contained in the given answer.

- The check is created/configured during config reading / init. / environment setup / in setupChild() / by ModuleFactory / by PackageObjectFactory.
- The check scans (the file contents of) the input source for tab character occurrences (which may lead to the creation of warning messages).
- The above happens at the command of TreeWalker (as it commences processing the input source). (mentioning beginTree is sufficient)
- The check is explicitly destroyed (by TreeWalker).

Task 3.2

Listed below are the eight calls that occur between TreeWalker and TabCharacterCheck during this specific scenario. Assign one point for every two correct calls (in the right order). Subtract one point for every two incorrect calls (and in case of one single incorrect call).

Note that in fact there is one more call, destroy(), but a bug in Extravis prevents this call from being shown in the MSV, even at 100visibility. Therefore, the reviewer should not take this call into consideration in either of the two subject groups.

Note that if the answer specifies nested calls within correct calls, they may be ignored by the reviewer as long as it is clear that they are nested.

contextualize

configure

init
getTokenNames
getDefaultTokens
setFileContents
beginTree
finishTree
destroy (not considered due to Extravis bug)

Task 3.3

visitToken
leaveToken

One correct: 2 points. Both correct: another point.

Subtract one point for each incorrect call. Also allow implementation-based reasoning.

Award 1 "bonus" point in case the motivation specifies that `IllegalInstantiationCheck` actually visits/checks for tokens (in the AST), whereas `TabCharacterCheck` does not (because it processes the file contents directly) – regardless of (in)correctness of the abovementioned identifiers.

Task 4.1

The following elements largely capture the essence of the error handling process. Assign one point for each element that is contained in the given answer. Note that mentioning the check's "mMessages"-field implies knowledge of both 2. and 3., and therefore yields points for both of these elements.

- A violation results in a call to `log()`, or in a read of `message.properties` for a human-readable format.
- The warning is internally stored/represented as an `api.LocalizedMessage`.
- The `LocalizedMessage` is added to the `api.LocalizedMessages` field (called "mMessages") of its Checker (in this case, `TreeWalker`) – note that there are multiple such repositories, not a global one!
- At the *end* of execution, the messages are relayed to the listeners (which each convert

the messages to different output types – human readable format / xml / etc.). Mentioning `fireErrors()` and its effect is also sufficient.

Task 4.2

The answer is no, which yields two points. The remaining points depend on the soundness of the motivation. The following are correct examples:

- Look for communication between the check and `api.LocalizedMessage(s)` in an execution trace.
- or: look for communication between the check and `api.DetailAST`.
- or: investigate the actual effect of `visitToken()`.
- or: find out what kind of human-readable message should result from violations in this check, and match this message with Checkstyle's (example) output. The true conclusion should be that there is no match: none of the warnings in this scenario's output relate to the check at hand.
- or: run and debug the application (e.g., using print statements, breakpoints, etc.).

If the answer is solely based on the interpretation of the `ws.notFollowed` file, the answer is partly correct because this trail will run cold – award two points in case the conclusion was "no", one point in case of a "yes".

Only two points are awarded if the reasoning is based on an understanding of what the check is looking for, and on the fact that the input source file contains no occurrences of whitespaces after tokens. No full score here because through this reasoning it cannot be determined that the check actually works.

Coupling measurements

Acquired through an automated analysis of the static call graph.

43 `TreeWalker`

35 `Checker`

8 `XMLLogger`

6 `CheckStyleTask`

```
6  DefaultLogger
6  DefaultConfiguration
-----
4  ConfigurationLoader$InternalLoader
3  PackageNamesLoader
2  ConfigurationLoader
2  DefaultContext
2  Main
1  PackageObjectFactory
1  PropertyCacheFile
```

Fanin/Fanout measurements

Acquired through an automated analysis of the static call graph.

```
api.AbstractFileSetCheck 32 11
api.AbstractLoader 8 0
api.AbstractViolationReporter 123 6
api.AuditEvent 25 6
api.AuditListener 5 0
api.CheckstyleException 11 0
api.Comment 20 0
api.Configurable 4 0
api.Configuration 10 0
api.Context 2 0
api.Contextualizable 3 0
api.DetailAST 375 0
api.FileContents 23 3
api.FileSetCheck 3 0
api.Filter 3 0
api.FilterSet 9 5
api.FullIdent 55 2
api.LocalizedMessage 14 0
api.LocalizedMessages 9 0
```

api.MessageDispatcher 14 0
api.Scope 11 0
api.ScopeUtils 36 3
api.SeverityLevel 14 0
api.SeverityLevelCounter 8 2
api.StrArrayConverter 1 0
api.TextBlock 19 0
api.TokenTypes 10 0
api.Utils 37 0
checks.AbstractFormatCheck 32 1
checks.AbstractOption 10 7
checks.AbstractOptionCheck 16 1
checks.AbstractTypeAwareCheck\$ClassAlias 4 4
checks.AbstractTypeAwareCheck\$ClassInfo 10 0
checks.AbstractTypeAwareCheck\$RegularClass 5 3
checks.AbstractTypeAwareCheck\$Token 10 3
checks.ArrayTypeStyleCheck 4 4
checks.BlockFrame 1 1
checks.blocks.AvoidNestedBlocksCheck 4 4
checks.blocks.BlockOption 1 1
checks.blocks.LeftCurlyOption 1 1
checks.blocks.NeedBracesCheck 4 3
checks.blocks.RightCurlyOption 1 1
checks.CheckUtils 15 7
checks.ClassFrame 1 1
checks.ClassResolver 2 0
checks.coding.AbstractIllegalCheck 4 0
checks.coding.AbstractSuperCheck\$MethodNode 4 0
checks.coding.AvoidInlineConditionalsCheck 6 4
checks.coding.DeclarationOrderCheck 5 5
checks.coding.DoubleCheckedLockingCheck 4 4
checks.coding.EmptyStatementCheck 4 3
checks.coding.FinalLocalVariableCheck 8 8
checks.coding.HiddenFieldCheck\$FieldFrame 6 0
checks.coding.IllegalThrowsCheck 6 6

checks.coding.InnerAssignmentCheck 4 4
checks.coding.ModifiedControlVariableCheck 8 7
checks.coding.MultipleStringLiteralsCheck\$stringInfo 3 0
checks.coding.MultipleVariableDeclarationsCheck 4 4
checks.coding.NestedIfDepthCheck 5 5
checks.coding.NestedTryDepthCheck 5 4
checks.coding.PackageDeclarationCheck 8 3
checks.coding.ParameterAssignmentCheck 8 6
checks.coding.ReturnCountCheck 9 7
checks.coding.SimplifyBooleanExpressionCheck 7 4
checks.coding.SimplifyBooleanReturnCheck 4 4
checks.coding.StringLiteralEqualityCheck 4 3
checks.coding.SuperCloneCheck 1 0
checks.coding.SuperFinalizeCheck 1 0
checks.DescendantTokenCheck 16 7
checks.design.FinalClassCheck\$classDesc 7 0
checks.design.InterfaceIsTypeCheck 6 4
checks.design.ThrowsCountCheck 6 6
checks.duplicates.ChecksumInfo 3 0
checks.duplicates.StrictDuplicateCodeCheck\$ChecksumGenerator 1 0
checks.duplicates.StrictDuplicateCodeCheck\$JavaChecksumGenerator 1 1
checks.duplicates.StrictDuplicateCodeCheck\$TextfileChecksumGenerator 2 1
checks.FileContentsHolder 6 1
checks.FrameStack 5 2
checks.GlobalFrame 1 1
checks.header.CrossLanguageRegexpHeaderCheck\$fileSetCheckViolationMonitor 2 1
checks.header.HeaderInfo 9 1
checks.header.HeaderViolationMonitor 2 0
checks.header.RegexpHeaderCheck\$CheckViolationMonitor 2 1
checks.header.RegexpHeaderInfo 5 2
checks.imports.AccessResult 1 0
checks.imports.AvoidStarImportCheck 4 4
checks.imports.Guard 3 0
checks.imports.PkgControl 4 2
checks.indentation.ExpressionHandler 129 41

checks.indentation.FinallyHandler 1 1
checks.indentation.HandlerFactory 4 1
checks.indentation.IndentLevel 24 0
checks.indentation.LineSet 9 0
checks.indentation.PrimordialHandler 6 2
checks.indentation.StaticInitHandler 1 1
checks.indentation.SwitchHandler 6 6
checks.indentation.TryHandler 3 3
checks.j2ee.AbstractInterfaceCheck 5 1
checks.j2ee.AbstractJ2eeCheck 17 5
checks.j2ee.BeanManagedMethodChecker 9 9
checks.j2ee.BeanMethodChecker 8 7
checks.j2ee.EntityBeanMethodChecker 16 10
checks.j2ee.HomeInterfaceMethodChecker 5 5
checks.j2ee.LocalInterfaceMethodChecker 2 2
checks.j2ee.MessageBeanMethodChecker 9 8
checks.j2ee.MethodChecker 32 16
checks.j2ee.PersistenceOption 1 1
checks.j2ee.RemoteInterfaceMethodChecker 2 2
checks.j2ee.SessionBeanMethodChecker 9 9
checks.j2ee.ThisParameterCheck 6 6
checks.j2ee.ThisReturnCheck 6 6
checks.j2ee.Utills 33 5
checks.javadoc.HtmlTag 7 0
checks.javadoc.JavadocMethodCheck\$ExceptionInfo 5 5
checks.javadoc.JavadocTag 14 0
checks.javadoc.Point 3 0
checks.LexicalFrame 6 0
checks.LineSeparatorOption 3 1
checks.MethodFrame 1 1
checks.metrics.AbstractComplexityCheck 13 4
checks.metrics.BooleanExpressionComplexityCheck 8 7
checks.metrics.ClassDataAbstractionCouplingCheck 3 2
checks.metrics.ClassFanOutComplexityCheck 3 1
checks.metrics.CyclomaticComplexityCheck 5 2

checks.metrics.JavaNCSSCheck 9 7
checks.metrics.JavaNCSSCheck\$Counter 2 0
checks.modifier.ModifierOrderCheck 4 3
checks.naming.ConstantNameCheck 4 4
checks.naming.LocalFinalVariableNameCheck 4 4
checks.naming.LocalVariableNameCheck 4 4
checks.naming.MethodNameCheck 3 1
checks.naming.ParameterNameCheck 4 2
checks.naming.StaticVariableNameCheck 4 4
checks.naming.TypeNameCheck 3 1
checks.sizes.AnonInnerLengthCheck 4 4
checks.sizes.ExecutableStatementCountCheck\$Context 4 0
checks.sizes.FileLengthCheck 4 2
CheckStyleTask\$Property 2 0
checks.UpperEllCheck 4 3
checks.whitespace.NoWhitespaceAfterCheck 5 5
checks.whitespace.OperatorWrapOption 1 1
checks.whitespace.PadOption 1 1
checks.whitespace.TabCharacterCheck 4 2
checks.whitespace.TypecastParenPadCheck 6 5
DefaultConfiguration 13 2
DefaultContext 5 0
DefaultLogger 8 6
doclets.CheckDocsDoclet 1 0
doclets.TokenTypesDoclet 0 0
filters.IntFilter 2 0
filters.IntMatchFilter 2 0
filters.IntRangeFilter 2 0
filters.SuppressionCommentFilter\$Tag 5 2
grammars.CommentListener 2 0
grammars.GeneratedJavaLexer 4 4
grammars.GeneratedJavaRecognizer 3 2
gui.AbstractTreeTableModel 6 2
gui.FileDrop 2 1
gui.FileDrop\$Listener 1 0

IEEE TRANSACTIONS ON SOFTWARE ENGINEERING, VOL. 0, NO. 0, JANUARY 2000

64

```
gui.JTreeTable$ListToTreeSelectionModelWrapper 3 0
gui.JTreeTable$TreeTableCellRenderer 3 1
gui.ParseTreeInfoPanel$FileDropListener 2 1
gui.ParseTreeInfoPanel$FileSelectionAction 1 1
gui.ParseTreeInfoPanel$ReloadAction 1 1
gui.ParseTreeModel 8 8
gui.TreeTableModel 7 0
ModuleFactory 2 0
PackageObjectFactory 3 1
PropertiesExpander 3 0
PropertyCacheFile 4 1
PropertyResolver 1 0
StringArrayReader 1 0
TreeWalker$SilentJavaRecognizer 1 1
```

March 8, 2010

DRAFT

TUD-SERG-2009-001
ISSN 1872-5392

