

Using Static Analysis and Static Measurement for Industrial Software Quality Evaluation

Luigi Lavazza

Dipartimento di Scienze Teoriche e Applicate
Università degli Studi dell'Insubria
Varese, Italy
Email: luigi.lavazza@uninsubria.it

Abstract—Business organizations that outsource software development need to evaluate the quality of the code delivered by suppliers. In this paper, we illustrate an experience in setting up and using a toolset for evaluating code quality for a company that outsources software development. The selected tools perform static code analysis and static measurement, and provide evidence of possible quality issues. To verify whether the issues reported by tools are associated to real problems, code inspections were carried out. The combination of automated analysis and inspections proved effective, in that several types of defects were identified. Based on our findings, the business company was able to learn what are the most frequent and dangerous types of defects that affect the acquired code: this knowledge is now being used on a regular basis to perform focused verification activities.

Index Terms—Software quality; Static analysis; Software measurement; Code clones; Code measures.

I. INTRODUCTION

Today, software is necessary for running practically any kind of business. However, poor quality software is generally expensive, because it can cause expensive failures, increased maintenance cost and security breaches. Hence, organizations that rely on software for running their business need to keep the quality of their software under control.

Many companies do not have the possibility or the will of developing the software they need. Hence, they outsource software development. In such case, an organization has no direct visibility and control of the development process; instead, they can only check the quality of the product they receive from developers.

In this paper, we report about an experience in setting up the toolset needed for evaluating the quality of the code provided by a supplier to an organization. This paper details and extends a previous report [1].

The software involved in the reported activities is used by the organization to run two Business-to-Consumer (B2C) portals. The organization needed to evaluate the quality of the supplied code; specifically, they wanted to check that the code was correctly structured, cleanly organized, well programmed, and free from defects that could cause failures or vulnerabilities that could be exploited by attackers. The organization had already in place a testing environment to evaluate the quality of the code from the point of view of

behavioral correctness. They wanted to complement test-based verification with checks on the internal qualities that could affect maintainability, fault-proneness and vulnerability.

To accomplish these goals, we selected a set of tools that provide useful indications based on static analysis and measurement of code. The toolset was intended to be used to evaluate two releases of the portal, and then to be set up at the company's premises, and used to evaluate the following releases.

The contributions of the paper are threefold:

- 1) We provide methodological guidance on the selection of a small set of tools that can provide quite useful insights on code quality.
- 2) We show how to use the tools and how to combine the automated analysis with manual verifications.
- 3) We provide some results, which can give the reader a measure of the results that can be achieved via the proposed approach.

Because of confidentiality constraints, in this paper we shall not identify the involved parties, and we shall omit some non-essential details.

The paper is structured as follows. In Section II, we provide some details concerning the evaluated software and the goals of the study. In Section III, the tools used for the static analysis and measurement are described. Section IV illustrates the methodological approach. In Section V, the results of the evaluation are described. Section VI provides some suggestions about the organization of a software development process that takes advantage of a static analysis and measurement toolset. Section VII illustrates the related work, while Section VIII draws some conclusions and outlines future work.

II. THE OBJECT OF THE EVALUATION

In this section, we describe the code evaluation issues that were tackled in the problem context.

The evaluation addressed two B2C portals, coded almost entirely in Java. The analyses concentrated exclusively on the Java code. Table I provides a few descriptive statistics concerning the two portals, where LOC is the total number of lines of code, while LLOC is the number of logical lines of code, i.e., the lines that contain executable code (excluding blank lines, comments, etc.).

TABLE I. CHARACTERISTICS OF THE ANALYZED PORTALS.

	Portal 1	Portal 2
Number of files	1,507	280
LLOC	100,375	37,467
LOC	202,249	55,934
Number of Classes	1,158	247
Number of Methods	13,351	5,370

The aim of the study consisted in evaluating the quality of the products, highlighting weaknesses and improvement opportunities. In this sense, it was important to spot the types of the most frequently recurring issues, rather than finding *all* the actual defects and issues.

It was also required that the toolset could be seamlessly transferred to the company's premises. To this end, open-source (or free to use) software was to be preferred.

Accordingly, we looked for tools that can

- Detect bad programming practices, based on the identification of specific code patterns.
- Detect bad programming practices, based on code measures (e.g., methods too long, classes excessively coupled, etc.).
- Detect duplicated code.
- Identify vulnerabilities.

III. TOOLS USED

After some search and evaluation, we selected the tools mentioned in Table II. These tools are described in some detail in the following sections.

A. Static Analysis for Identifying Defects

SpotBugs (formerly known as FindBugs) is a program which uses static analysis to look for bugs in Java code [2][3]. SpotBugs looks for code patterns that are most likely associated to defects. For instance, SpotBugs is able to identify the usage of a reference that is possibly null.

SpotBugs was chosen because it is open-source and one among the best known tools of its kind. Besides, SpotBugs proved to be quite efficient: on a reasonably powerful laptop, it took less than a minute to analyze Portal 1.

SpotBugs provides "confidence" and "rank" for each of the issued warnings. Confidence indicates how much SpotBugs is confident that the issued warning is associated to a real problem; confidence is expressed in a three-level ordinal scale (high, medium, low). The rank indicates how serious the problem is believed to be. The rank ranges from 1 (highest) to 20 (lowest); SpotBugs also indicates levels: "scariest" ($1 \leq \text{rank} \leq 4$), "scary" ($5 \leq \text{rank} \leq 9$), "worrying" ($10 \leq \text{rank} \leq 14$), "of concern" ($15 \leq \text{rank} \leq 20$).

SpotBugs is quite expressive. Figure 1 shows an example of the output provided by SpotBugs. The upper-left window shows the hierarchy of bugs, organized by type and/or by rank. The upper-right window shows the code, highlighting the part that is responsible of the bug. The bottom right window gives the "conceptual" description of the bug type: in the

show case the problem is due to using `String.equals` to compare objects having different types. The bottom-left windows explains why the shown code was recognized as responsible for the bug: it reports the file and line number where the problem was found, and which code elements are involved in the detected issue.

The provided information indicates clearly which problem was identified and where: usually it is easy for developers to manually check whether the detected bug is a real problem or a false positive. In fact, static analysis cannot provide exclusively correct indications. To avoid overlooking important problems whose occurrence the tool cannot decide, the provided results are usually somewhat "pessimistic." That is, some indications are false positives, and do not correspond to actual bugs. Because of this intrinsic characteristic of static analysis tools, manual verification of the reported issues is advisable. To get an idea of the level of precision yielded by SpotBugs, experimental evaluations found that the SpotBugs' precision is between 58% and 80% [4].

SpotBugs can save results in XML files. It is thus possible to write specific code to read and analyze such results. As an example of this kind of analysis, we wrote a script in the R language [5] to graphically show the distribution of bugs through code files. The result is given in Figure 2, where the file names have been hidden for confidentiality. The size of each box is proportional to the maximum gravity of the bugs it contains, while the color indicates the number of bugs: the more bugs in the file, the darker the box. This kind of representation was useful to drive the attention of project managers onto the most critical files.

B. Static Analysis for Identifying Vulnerabilities

Having selected SpotBugs as a static analyzer, it was fairly natural to equip it with FindSecBugs [6], a plug-in that addresses security problems.

FindSecBugs works like FindBugs and SpotBugs, looking for code patterns that can be associated to security issues, with reference to problems reported by the Open Web Application Security Project (OWASP) [7] and the weaknesses listed in the Common Weaknesses Enumeration (CWE) [8].

C. Static Measures of Code

Several tools are available to measure the most relevant characteristics of code, including size, complexity, coupling, cohesion, etc.

SourceMeter [9] was chosen because it is free to use, efficient and provides many measures, including all the most popular and relevant. Of the many measures that SourceMeter can provide, we used those listed in Table III. Specifically, we selected a relatively small set of measures that cover all the potentially interesting aspects of code, including size, complexity, cohesion, and coupling.

TABLE II. TOOLS USED AND THEIR PURPOSE.

Purpose	Tool	Main features
Identify defects	SpotBugs	Static analysis is used to identify code patterns that are likely associated to defects.
Collect static measures	SourceMeter	Static measurement is applied at different granularity levels (class, method, etc.) to provide a variety of measures.
Detect code clones	SourceMeter	Structurally similar code blocks are identified.
Identify security issues	FindSecBugs	A plug-in for FindBugs, specifically oriented to identifying vulnerable code.

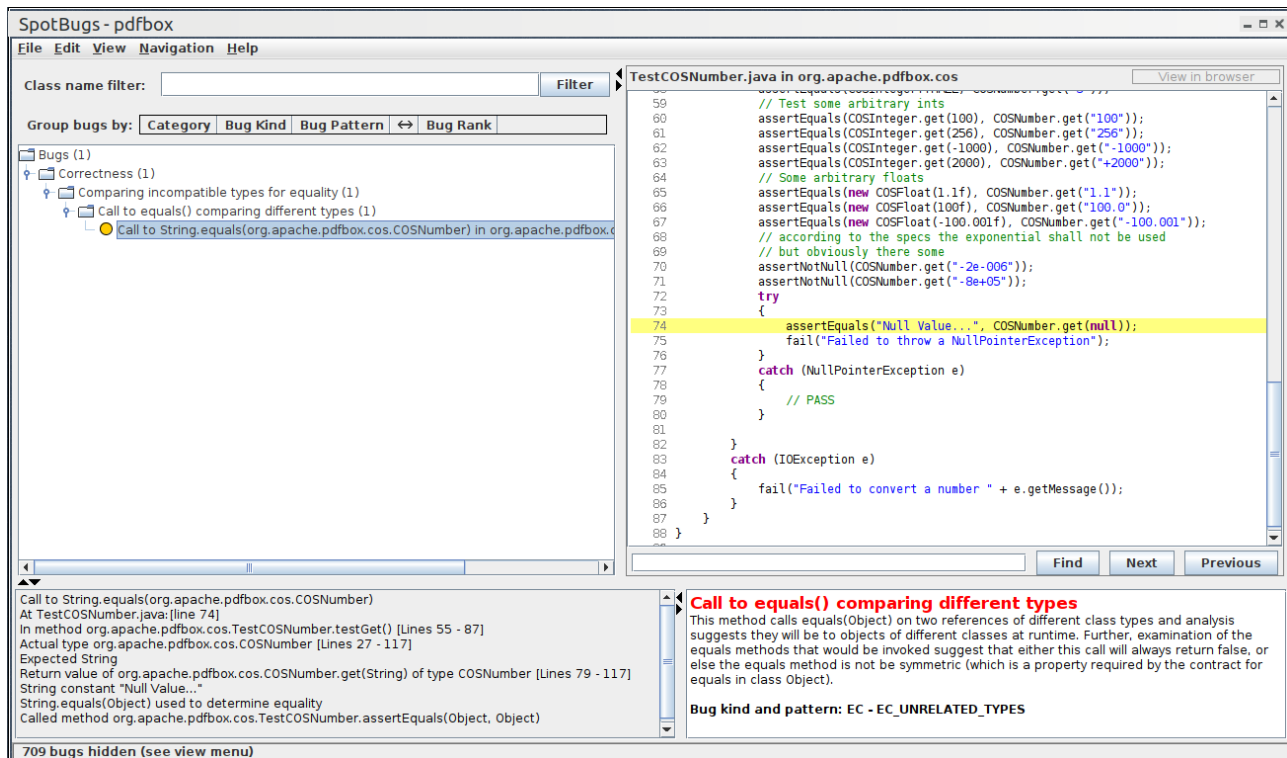


Fig. 1. A screenshot of SpotBugs.

TABLE III. Measures provided by SourceMeter and used in the evaluation.

Metric name	Abbreviation	Class	Method
Lack of Cohesion in Methods	LCOM5	X	
McCabe's Cyclomatic Complexity	McCC		X
Nesting Level Else-If	NLE	X	X
Weighted Methods per Class	WMC	X	
Coupling Between Object classes	CBO	X	
Response set For Class	RFC	X	
Depth of Inheritance Tree	DIT	X	
Number of Children	NOC	X	
Logical Lines of Code	LLOC	X	X
Number of Attributes	NA	X	
Number of Parameters	NUMPAR		X

D. Code Clone Detection

Noticeably, SourceMeter is also able to detect code clones. Specifically, SourceMeter is capable of identifying the so-called Type-2 clones, i.e. code fragments that are structurally identical, but may differ in variable names, literals, identifiers, etc.

IV. THE METHOD

Since the most interesting properties of code are undecidable, tools that perform static analysis often issue warnings concerning problems that are likely—but not certain—to occur. In practice, the issues reported by static analysis tools can be false positives. Therefore, we always inspected manually the code that had been flagged as possibly incorrect by the tools.

The inspection of issues reported by static analysis tools was generally quick and required little effort, because the tools provide a very clear indication of the nature of the problem and where it is. Sometimes the inspections required a substantial amount of effort, but that is probably because we were not at all familiar with the code: the developers of the code would have completed the inspection much more quickly.

Also code that is characterized by unusual values of static measures needed inspections. For instance, consider a method that has unusually high McCabe complexity: only via manual inspection we can check whether the program was badly structured or the coded algorithm is intrinsically complex.

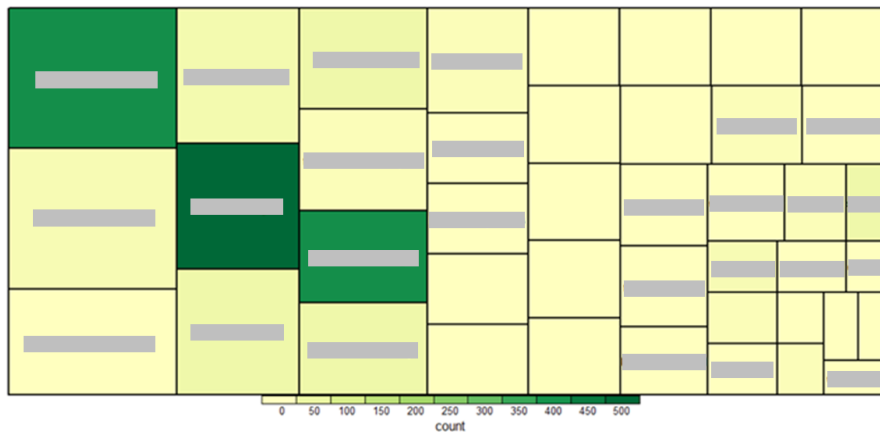


Fig. 2. A graphical representation of the distribution of bugs through code files.

In several cases static measures indicated as suspicious code elements that had already been flagged by other measures or by static analysis: in those cases, recognizing the nature of the problem (if any) was practically immediate.

Problem detection was performed as described in Figure 3. The real problems identified via the process described in Figure 3 were classified according to their type, so that the company that asked for the code quality analysis could focus improvement efforts on the most frequent and serious problems.

V. RESULTS

Here, we describe the code quality problems that were identified.

A. Warnings issued by SpotBugs

Tables IV and V illustrate the number of warnings that SpotBugs issued for the analyzed code, respectively by confidence level and by rank. In Table IV, the density indicates the number of warnings per line of code.

SpotBugs also classifies warning by type (for additional information on warning types, see [10]). Table VI illustrates the warnings we obtained, by type. It can be noticed that most warnings concerned security (types “Security” and “Malicious code vulnerability”).

1) *Results deriving from the inspection of SpotBugs warnings:* The effort allocated to the project did not allow analyzing all the warnings issued by SpotBugs. Therefore, we inspected the code where SpotBugs had identified issues ranked “scary” and “scariest.” Specifically, we analyzed the warnings described in Table VII.

Our inspections revealed several code quality problems:

- The existence of problems matching the types of warning issued by SpotBugs was confirmed.
- Some language constructs were not used properly. For instance, class `Boolean` was incorrectly used instead of `boolean`; objects of type `String` were used instead of `boolean` and `enumeration`; etc.

- We found redundant code, i.e., some pieces of code were unnecessarily repeated, even where avoiding code duplication—e.g., via inheritance or even simply by creating methods that could be used in different places—would have been easy to use and definitely convenient.
- We found some pieces of code that were conceptually incorrect. The types of defect were not of any type that a static analyzer could find, but were quite apparent when inspecting the code.

Concerning the correctness of warnings issued by SpotBugs, we found just one false positive: the “comparison of `String` objects using `==` or `!=`” was not an error, in the examined case. We also found that the four instances of “Method ignores return value” were of little practical consequences. In summary, the great majority of warnings indicated real problems, which could cause possibly serious consequences. The remaining warning indicated situations where a better coding discipline could make the code less error prone, if applied systematically.

2) *Results deriving from the inspection of FindSecBugs warnings:* The great majority of the security warnings (types “Security” and “Malicious code vulnerability”) were ranked by FindSecBugs as not very worrying. Specifically, no “scariest” warning was issued, and only one “scary” warning was issued. Therefore, we inspected the only “scary” warning (rank 7, see Table VIII), and all the warnings at the highest rank of the level “troubling” (rank 10, see Table VIII).

We found that all the warnings pointed to code that had security problems. In many cases, SpotBugs documentation provided quite straightforward ways for correcting the code.

B. Inspection of code elements having measures beyond threshold

Static measures concerning size, complexity, cohesion, coupling, among others, are expected to provide indications on the quality of code. In fact, one expects that code characterized by large size, high complexity, low cohesion, strong coupling and similar “bad” characteristics is error-prone. Accordingly, we inspected code elements having measures definitely out of

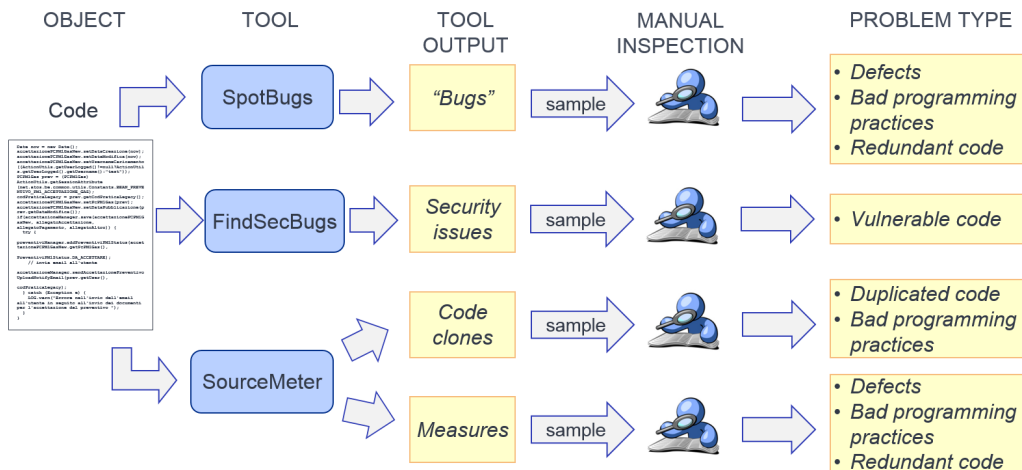


Fig. 3. The evaluation process: problem detection phase.

TABLE IV. SPOTBUGS WARNINGS BY CONFIDENCE.

Metric	Portal 1		Portal 2	
	Warnings	Density	Warnings	Density
High Confidence	68	0.07%	50	0.13%
Medium Confidence	774	0.77%	502	1.34%
Low Confidence	824	0.82%	420	1.12%
Total	1666	1.66%	972	2.59%

TABLE V. SPOTBUGS WARNINGS BY RANK

Rank	1	6	7	8	9	10	11	12	14	15	16	17	18	19	20
Portal 1	24	9	1	42		12	39	21	2	604	17	129	562	82	122
Portal 2	10			6	1	7	8	27	18	191	10	59	401	66	168

TABLE VI. SPOTBUGS WARNINGS BY TYPE.

Warning Type	Portal 1		Portal 2	
	Number	Percentage	Number	Percentage
Bad practice	73	4.38%	81	8.33%
Correctness	91	5.46%	30	3.09%
Experimental	0	0.00%	1	0.10%
Internationalization	16	0.96%	39	4.01%
Malicious code vulnerability	496	29.77%	316	32.51%
Multithreaded correctness	28	1.68%	1	0.10%
Performance	74	4.44%	143	14.71%
Security	631	37.88%	215	22.12%
Dodgy code	257	15.43%	146	15.02%
Total	1666	100%	972	100%

the usually considered safe ranges. Specifically, we considered McCabe complexity [11], Logical Lines of Code and Response for Class (RFC) [12] as possibly correlated with problems. In fact, we also looked at Coupling Between Objects, Lack of Cohesion in Methods and Weighted Method Count, but these measures turned out to provide no additional information with respect to the aforementioned three measures, i.e., they pointed to the same classes or methods identified as possibly problematic by the aforementioned measures.

We found that several methods featured McCabe complexity well over the threshold that is generally considered safe. Figure 4 shows the distributions of McCabe complexity of methods (excluding setters and getters) together with two thresholds representing values that should and must not be exceeded, according to common knowledge [11][13][14][15]. Specifically, we found methods with McCabe complexity close to 200.

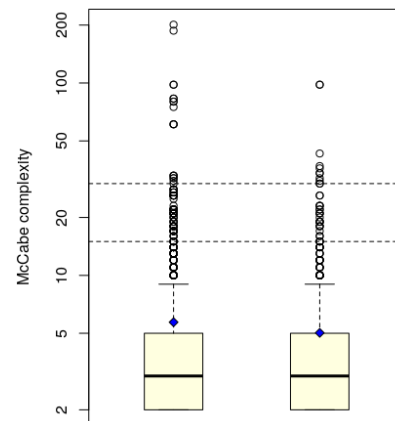


Fig. 4. Boxplots illustrating the distributions of McCabe complexity in the two portals (blue diamonds indicate mean values). The scale is logarithmic.

TABLE VII. SPOTBUGS WARNINGS THAT WERE VERIFIED MANUALLY.

Rank	Type	Occurrences	
		Portal 1	Portal 2
1	Suspicious reference comparison	10	9
1	Call to equals() comparing different types	14	1
6	Possible null pointer dereference	8	-
8	Possible null pointer dereference	-	2
8	Method ignores return value	-	4
9	Comparison of String objects using == or !=	-	1

TABLE VIII. FINDSECBUGS WARNINGS THAT WERE VERIFIED MANUALLY.

Rank	Type	Occurrences	
		Portal 1	Portal 2
7	HTTP response splitting vulnerability	1	-
10	Cipher with no data integrity	4	2
10	ECB mode is insecure	4	2
10	URL Connection Server-Side Request Forgery and File Disclosure	1	-
10	Unvalidated Redirect	2	-
10	Request Dispatcher File Disclosure	-	1

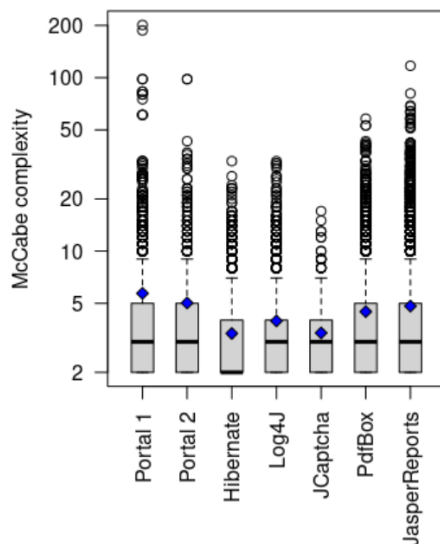


Fig. 5. Boxplots illustrating the distributions of McCabe complexity in the two portals in comparison with a few popular open-source products. The scale is logarithmic.

The values of McCabe complexity in the evaluated portals may seem exceedingly large. To assess these values, we compared them to the McCabe complexity of a few popular open-source projects, which are widely used and generally considered of good quality. The comparison is given in Figure 5. It can be seen that the distributions of McCabe complexity values in the evaluated portals is not very different from other programs'. However, in all cases the methods having McCabe complexity greater than 10 are considered outliers: they deserve careful quality analysis.

When considering size, we found several classes featuring over 1000 LLOC; the largest class contained slightly less than 6000 LLOC. When considering RFC, we found 12 classes having RFC greater than 200. Interestingly, the class with the highest RFC (709) was also the one containing the method with the greatest McCabe complexity. The biggest class con-

tained the second most complex method. These results were not surprising, since it is known that several measures are correlated to size.

Inspections revealed that the classes and methods featuring excessively high values of LLOC, RFC and McCabe complexity were all affected by the same problem. The considered code had to deal with several types of services, which were very similar under several respects, although each one had its own specificity. The analyzed code ignored the similarities among the services to be managed, so that the code dealing with similar service aspects was duplicated in multiple methods. The code could have been organized differently using basic object-oriented features: a generic class could collect the features that are common to similar services, and a specialized class for every service type could take care of the specificity of different service types.

In conclusion, by inspecting code featuring unusual static measures, we found design problems, namely inheritance and late binding were not used where it was possible and convenient.

C. Inspection of duplicated code

SourceMeter was also used to find duplicated code. Specifically, structurally similar blocks of 10 or more lines of code were looked for. Many duplicated blocks were found. For instance, in Portal 1, 434 duplicated blocks were found. In many cases, blocks included more than one hundred lines. The largest duplicated blocks contained 205 lines. A small minority of detections concerned false positives.

We found three types of duplications:

- Duplicates within the same file. That is, the same code was found in different parts of the same file (or the same class, often).
- Duplicates in different files. That is, the same code fragment was found in different files (of the same portal).
- Duplicates in different portals. That is, the same code fragment was found in files belonging to different portal.

Duplicates of type c) highlighted the existence of versioning problems: different versions of the same class were used in the two portals.

Duplicates of types a) and b) pointed to the same type of problem already identified, i.e., not using inheritance to factor code that can be shared among classes dealing with similar services. Concerning this issue, it is worth noting that static measures revealed a general problem with the design of code,

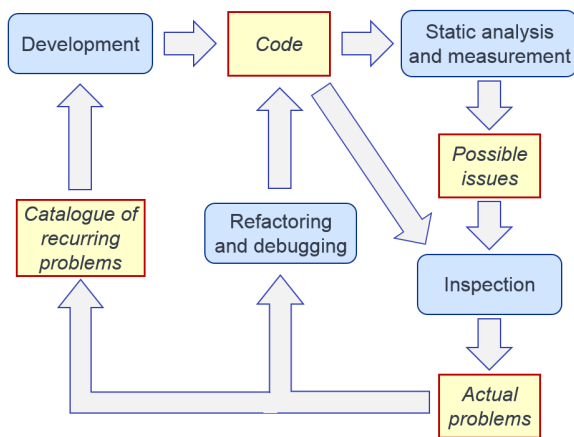


Fig. 6. Suggested Development Process.

but were not able to indicate precisely which parts of the code could be factorized. On the contrary, duplicated code detection was quite effective in identifying all the cases where code could be factorized, with little need of inspecting the code. In this sense, code clone detection added some value to inspections aiming at understanding the reasons for ‘out of range’ measures.

VI. SUGGESTIONS FOR IMPROVING THE DEVELOPMENT PROCESS

Given the results described in Section V, it seems convenient that the capabilities of static analysis and measurement tools are exploited on a regular basis. To this end, we can observe that two not exclusive approaches are possible.

1) *Evaluation of code:* The toolset can be used to evaluate the released code as described in Section IV. However, it would be advisable that developers verify their own code via SpotBugs and SourceMeter *before* releasing it: in such a way, a not negligible number of bugs would be removed even before testing and other Verification&Validation activities, thus saving time and effort. With respect to the evaluation described in Section IV, where just a sample of the issues reported by the tools were inspected, in the actual development process all issues should be inspected.

2) *Prevention:* The practice of issue identification and verification leads to identifying the most frequently recurring types of problems. It is therefore possible to compile a catalogue of the most frequent and dangerous problems. Accordingly, programmers could be instructed to carefully avoid such issues. This could imply teaching programmers specific techniques and good programming practices.

As a result of the considerations illustrated above, software development activities could be organized as described in Figure 6. If development is outsourced, as in the cases described in this paper, the catalogue of recurrent problems could be used as part of the contract annex that specifies the required code quality level.

Finally, it is worth noting that the proposed approach can be applied in practically any type of lifecycle. For instance,

in an agile development environment, the proposed evaluation practices could be applied at the end of every sprint.

VII. RELATED WORK

The effectiveness of using automated static analysis tools for detecting and removing bugs was documented by Zheng et al. [16]. Among other facts, they found that the cost per detected fault is of the same order of magnitude for static analysis tools and inspections, and the defect removal yield of static analysis tools is not significantly different from that of inspections.

Thung et al. performed an empirical study to evaluate to what extent could field defects be detected by FindBugs and similar tools [17]. To this end, FindBugs was applied to three open-source programs (Lucene, Rhino and AspectJ). The study by Thung et al. takes into consideration only known bugs, and is performed on open-source programs. On the contrary, we analyzed programs developed in an industrial context, and relied on manual inspection to identify actual bugs.

Habib and Pradel performed a study to determine how many of all real-world bugs do static bug detectors find [18]. They used three static bug detectors, including SpotBugs, to analyze a version of the Defects4J dataset that consisted of 15 Java projects with 594 known bugs. They found that static bug detectors find a small but non-negligible amount of all bugs.

Vetrò et al. [19] evaluated the accuracy of FindBugs. The code base used for the evaluation consisted of Java projects developed by students in the context of an object-oriented programming course. The code is equipped with acceptance tests written by teachers of the course in such a way that all functionalities are checked. To determine true positives, they used temporal and spatial coincidence: an issue was considered related to a bug when an issue disappeared at the same time as a bug get fixed (according to tests). In a later paper [20] Vetrò et al. repeated the analysis, with a larger code set and performing inspections concerning four types of issues found by FindBugs, namely the types of findings that are considered more reliable.

Tomassi [21] considered 320 Java bugs from the BugSwarm dataset, and determine which of these bugs can potentially be found by SpotBugs and another analyzer—namely, ErrorProne (<https://github.com/google/error-prone>)— and how many are indeed detected. He found that 40.3% of the bugs were of types that SpotBugs should detect, but only one of such bugs was actually detected by SpotBugs.

In general, the papers mentioned above have goals and use methods that are somewhat similar to ours, but are nonetheless different in important respects. A work that shares context, goals and methods with ours was reported by Steidl et al. [22]. They observed that companies often use static analyses tools, but they do not learn from results, so that they fail to improve code quality. Steidl et al. propose a continuous quality control process that combines measures, manual action, and a close cooperation between quality engineers, developers, and managers. Although there are evident differences between the work by Steidl et al. and the work reported in this paper

(for instance, the situation addressed by Steidl et al. does not involve outsourcing), the suggestions for improving the development process given in Section VI are conceptually coherent with the proposal by Steidl et al.

Similarly, Wagner et al. [23] performed an evaluation of the effectiveness of static analysis tools in combination with other techniques (including testing and reviews). They observed that a combination of the usage of bug finding tools together with reviews and tests is advisable if the number of false positives is low, as in fact is in the cases we analyzed (many false positives would imply that a relevant effort is wasted).

An alternative to static analyzers like SpotBugs is given by tools that detect the presence of “code smells” [24] in code. A comparison of these types of tools was performed by applying SpotBugs and JDeodorant [25][26] to a set of set of open-source applications [4]. The study showed that the considered tools can help software practitioners detect and remove defects in an effective way, to limit the amount of resources that would otherwise be spent in more cost-intensive activities, such as software inspections. Specifically, SpotBugs appeared to detect defects with good Precision, hence manual inspection of the code flagged defective by SpotBugs becomes cost-effective.

Another empirical study evaluated the persistence of SpotBugs issues in open-source software evolution [27]. This study showed that around half the issues discovered by SpotBugs are actually removed from code. This fact is interpreted as a confirmation that SpotBugs identifies situations that are considered worth correcting by developers.

VIII. CONCLUSIONS

Evaluating the quality of software is important in general, and especially for business organization that outsource development, and do not have visibility and control of the development process. Software testing can provide some kind of quality evaluations, but to a limited extent. In fact, some aspects of code quality (e.g., whether the code is organized in a way that favors maintainability) cannot be assessed via testing.

This paper describes an approach to software quality evaluation that consists of two phases: in the first phase, tools are used to identify possible issues in the code; in the second phase, code is manually inspected to verify whether the reported issues are associated to real problems. The tools used are of two kinds: the first performs static analysis of code looking for patterns that are likely associated to problematic code; the second type yields measures of static code properties (like size, complexity, cohesion, coupling etc.), thus helping identifying software elements having excessive, hence probably problematic, characteristics.

The mentioned approach was applied to the code of the web portals used by a European company to let its customers use a set of services. The experience was successful, as tool-driven inspections uncovered several types of defects. In the process, the tools (namely SpotBugs and SourceMeter) identified problems of inherently different nature, hence it is advisable to use both types of tools.

Based on our findings, the business company was able to learn what are the most frequent and dangerous types of defects that affect the acquired code: this knowledge is being used to perform focused verification activities.

The proposed approach and toolset (possibly composed of equivalent tools) can be useful in several contexts where code quality evaluation is needed. Noticeably, the proposed approach can be used in different types of development process, including agile processes.

Among the future possible evolutions of this work, the most intriguing one concerns studying the possibility of replacing inspection via some sort of AI-based models that can discriminate false positives and true problems. Specifically, collecting the results of manual inspections driven by SpotBugs could lead to some sort of supervised learning.

Another interesting evolution of the presented work involves using automatic code smell detectors, like JDeodorant [28], for instance, in addition to SpotBugs.

ACKNOWLEDGMENT

This work has been partially supported by the “Fondo di ricerca d’Ateneo” of the Università degli Studi dell’Insubria.

REFERENCES

- [1] L. Lavazza, “Software quality evaluation via static analysis and static measurement: an industrial experience,” in *The Fifteenth International Conference on Software Engineering Advances—ICSEA*, 2020, pp. 55–60.
- [2] D. Hovemeyer and W. Pugh, “Finding bugs is easy,” *ACM Sigplan notices*, vol. 39, no. 12, 2004, pp. 92–106.
- [3] <https://spotbugs.github.io/> [retrieved: August 2020].
- [4] L. Lavazza, S. Morasca, and D. Tosi, “Comparing static analysis and code smells as defect predictors: an empirical study,” in *17th IFIP International Conference on Open Source Systems (OSS)*. Springer International Publishing, 2021, pp. 1–15.
- [5] R. C. Team et al., “R: A language and environment for statistical computing,” 2013.
- [6] <https://find-sec-bugs.github.io/> [retrieved: August 2020].
- [7] “Open Web Application Security Project® (OWASP),” <https://www.owasp.org> [retrieved: August 2020].
- [8] “Common Weaknesses Enumeration,” <https://cwe.mitre.org> [retrieved: August 2020].
- [9] R. Ferenc, L. Langó, I. Siket, T. Gyimóthy, and T. Bakota, “Source meter sonar qube plug-in,” in *2014 IEEE 14th International Working Conference on Source Code Analysis and Manipulation*, 2014, pp. 77–82.
- [10] <https://spotbugs.readthedocs.io/en/stable/bugDescriptions.html> [retrieved: August 2020].
- [11] T. J. McCabe, “A complexity measure,” *IEEE Transactions on software Engineering*, no. 4, 1976, pp. 308–320.
- [12] S. R. Chidamber and C. F. Kemerer, “A metrics suite for object oriented design,” *IEEE Transactions on software engineering*, vol. 20, no. 6, 1994, pp. 476–493.
- [13] A. H. Watson, D. R. Wallace, and T. J. McCabe, *Structured testing: A testing methodology using the cyclomatic complexity metric*. NIST – National Institute of Standards and Technology, 1996, vol. 500, no. 235.
- [14] M. Bray, K. Brune, D. A. Fisher, J. Foreman, and M. Gerken, “C4 software technology reference guide—a prototype.” Carnegie-Mellon Univ., Pittsburgh PA, Software Engineering Inst, Tech. Rep., 1997.
- [15] JPL, “JPL Institutional Coding Standard for the C Programming Language,” Jet Propulsion Laboratory, Tech. Rep., 2009.
- [16] J. Zheng, L. Williams, N. Nagappan, W. Snipes, J. P. Hudspohl, and M. A. Vouk, “On the value of static analysis for fault detection in software,” *IEEE transactions on software engineering*, vol. 32, no. 4, 2006, pp. 240–253.

- [17] F. Thung, L. Lucia, D. Lo, L. Jiang, P. Devanbu, and F. Rahman, "To what extent could we detect field defects? an extended empirical study of false negatives in static bug-finding tools," *Automated Software Engineering*, vol. 22, no. 4, 2015, pp. 561–602.
- [18] A. Habib and M. Pradel, "How many of all bugs do we find? a study of static bug detectors," in *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*, 2018, pp. 317–328.
- [19] A. Vetrò, M. Torchiano, and M. Morisio, "Assessing the precision of FindBugs by mining java projects developed at a university," in *2010 7th IEEE Working Conference on Mining Software Repositories (MSR 2010)*. IEEE, 2010, pp. 110–113.
- [20] A. Vetrò, M. Morisio, and M. Torchiano, "An empirical validation of findbugs issues related to defects," in *15th Annual Conference on Evaluation and Assessment in Software Engineering (EASE 2011)*. IET, 2011, pp. 144–153.
- [21] D. A. Tomassi, "Bugs in the wild: examining the effectiveness of static analyzers at finding real-world bugs," in *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2018, pp. 980–982.
- [22] D. Steidl, F. Deissenboeck, M. Poehlmann, R. Heinke, and B. Uhin-Mergenthaler, "Continuous software quality control in practice," in *2014 IEEE International Conference on Software Maintenance and Evolution*. IEEE, 2014, pp. 561–564.
- [23] S. Wagner, J. Jürjens, C. Koller, and P. Trischberger, "Comparing bug finding tools with reviews and tests," in *IFIP International Conference on Testing of Communicating Systems*. Springer, 2005, pp. 40–55.
- [24] M. Fowler, *Refactoring: improving the design of existing code*. Addison-Wesley Professional, 2018.
- [25] M. Fokaefs, N. Tsantalis, and A. Chatzigeorgiou, "Jdeodorant: Identification and removal of feature envy bad smells," in *2007 IEEE International Conference on Software Maintenance*. IEEE, 2007, pp. 519–520.
- [26] "JDeodorant website," 2020. [Online]. Available: <https://github.com/tsantalis/JDeodorant>
- [27] L. Lavazza, D. Tosi, and S. Morasca, "An empirical study on the persistence of spotbugs issues in open-source software evolution," in *13th International Conference on the Quality of Information and Communications Technology – QUATIC*, 2020.
- [28] N. Tsantalis, T. Chaikalis, and A. Chatzigeorgiou, "Ten years of jdeodorant: Lessons learned from the hunt for smells," in *2018 IEEE 25th International Conference on Software Analysis, Evolution and Reengineering (SANER)*. IEEE, 2018, pp. 4–14.