

A Systematic Literature Review on Fault Prediction Performance in Software Engineering

Tracy Hall, Sarah Beecham, David Bowes, David Gray, and Steve Counsell

Abstract—*Background:* The accurate prediction of where faults are likely to occur in code can help direct test effort, reduce costs, and improve the quality of software. *Objective:* We investigate how the context of models, the independent variables used, and the modeling techniques applied influence the performance of fault prediction models. *Method:* We used a systematic literature review to identify 208 fault prediction studies published from January 2000 to December 2010. We synthesize the quantitative and qualitative results of 36 studies which report sufficient contextual and methodological information according to the criteria we develop and apply. *Results:* The models that perform well tend to be based on simple modeling techniques such as Naive Bayes or Logistic Regression. Combinations of independent variables have been used by models that perform well. Feature selection has been applied to these combinations when models are performing particularly well. *Conclusion:* The methodology used to build models seems to be influential to predictive performance. Although there are a set of fault prediction studies in which confidence is possible, more studies are needed that use a reliable methodology and which report their context, methodology, and performance comprehensively.

Index Terms—Systematic literature review, software fault prediction

1 INTRODUCTION

THIS Systematic Literature Review (SLR) aims to identify and analyze the models used to predict faults in source code in 208 studies published between January 2000 and December 2010. Our analysis investigates how model performance is affected by the context in which the model was developed, the independent variables used in the model, and the technique on which the model was built. Our results enable researchers to develop prediction models based on best knowledge and practice across many previous studies. Our results also help practitioners to make effective decisions on prediction models most suited to their context.

Fault¹ prediction modeling is an important area of research and the subject of many previous studies. These studies typically produce fault prediction models which allow software engineers to focus development activities on fault-prone code, thereby improving software quality and

making better use of resources. The many fault prediction models published are complex and disparate and no up-to-date comprehensive picture of the current state of fault prediction exists. Two previous reviews of the area have been performed in [1] and [2].² Our review differs from these reviews in the following ways:

- *Timeframes.* Our review is the most contemporary because it includes studies published from 2000–2010. Fenton and Neil conducted a critical review of software fault prediction research up to 1999 [1]. Catal and Diri's [2] review covers work published between 1990 and 2007.
- *Systematic approach.* We follow Kitchenham and Charters [3] original and rigorous procedures for conducting systematic reviews. Catal and Diri did not report on how they sourced their studies, stating that they adapted Jørgensen and Shepperd's [4] methodology. Fenton and Neil did not apply the systematic approach introduced by Kitchenham and Charters [3] as their study was published well before these guidelines were produced.
- *Comprehensiveness.* We do not rely on search engines alone and, unlike Catal and Diri, we read through relevant journals and conferences paper-by-paper. As a result, we analyzed many more papers.
- *Analysis.* We provide a more detailed analysis of each paper. Catal and Diri focused on the context of studies, including: where papers were published, year of publication, types of metrics used, datasets used, and modeling approach. In addition, we report

1. The term "fault" is used interchangeably in this study with the terms "defect" or "bug" to mean a static fault in software code. It does not denote a "failure" (i.e., the possible result of a fault occurrence).

- T. Hall and S. Counsell are with the Department of Information Systems and Computing, Brunel University, Uxbridge, Middlesex UB8 3PH, United Kingdom. E-mail: {tracy.hall, steve.counsell}@brunel.ac.uk.
- S. Beecham is with Lero—The Irish Software Engineering Research Centre, University of Limerick, Tierney Building, Limerick, Ireland. E-mail: sarah.beecham@lero.ie.
- D. Bowes and D. Gray are with the Science and Technology Research Institute, University of Hertfordshire, Hatfield, Hertfordshire AL10 9AB, United Kingdom. E-mail: {d.h.bowes, d.gray}@herts.ac.uk.

Manuscript received 21 Oct. 2010; revised 10 July 2011; accepted 13 Sept. 2011; published online 30 Sept. 2011.

Recommended for acceptance by D. Sjøberg.

For information on obtaining reprints of this article, please send e-mail to: tse@computer.org, and reference IEEECS Log Number TSE-2010-10-0312. Digital Object Identifier no. 10.1109/TSE.2011.103.

2. Note that two referencing styles are used throughout this paper; [ref#] refers to papers in the main reference list while [Sref#] refers to papers in the separate systematic literature review list, located before the main reference list.

TABLE 1
The Research Questions Addressed

Research Questions		Motivation
RQ1	How does context affect fault prediction?	Context has been shown to be a key factor in the comparative use of software metrics in general [5]. Context is important in fault prediction modelling as it can affect the performance of models in a particular context and the transferability of models between contexts. Currently the impact context variables have on the transferability of models is not clear. This makes it difficult for potential model users to select models that will perform well in a particular context. We aim to present a synthesis of current knowledge on the impact of context on models and the transferability of models.
RQ2	Which independent variables should be included in fault prediction models?	There are a range of independent variables that have been used in fault prediction models. Currently the impact individual independent variables have on model performance is not clear. Although the performance of independent variables has been investigated within individual studies, no comparison of performance across studies has been done. This makes it difficult for model builders to make informed decisions about the independent variables on which to base their models. We aim to present a synthesis of current knowledge on the impact independent variables have on models.
RQ3	Which modeling techniques perform best when used in fault prediction?	Fault prediction models are based on a wide variety of both machine learning and regression modelling techniques. Currently the impact modelling technique has on model performance is not clear. Again, the performance of modelling techniques has been investigated within individual studies, but no comparison of performance across studies has been done. This makes it difficult for model builders to make effective technique selections. We aim to present a synthesis of current knowledge on the impact of modelling technique on model performance.

on the performance of models and synthesize the findings of studies.

We make four significant contributions by presenting:

1. A set of 208 studies addressing fault prediction in software engineering from January 2000 to December 2010. Researchers can use these studies as the basis of future investigations into fault prediction.
2. A subset of 36 fault prediction studies which report sufficient contextual and methodological detail to enable these studies to be reliably analyzed by other researchers and evaluated by model users planning to select an appropriate model for their context.
3. A set of criteria to assess that sufficient contextual and methodological detail is reported in fault prediction studies. We have used these criteria to identify the 36 studies mentioned above. They can also be used to guide other researchers to build credible new models that are understandable, usable, replicable, and in which researchers and users can have a basic level of confidence. These criteria could also be used to guide journal and conference reviewers in determining that a fault prediction paper has adequately reported a study.
4. A synthesis of the current state of the art in software fault prediction as reported in the 36 studies satisfying our assessment criteria. This synthesis is based on extracting and combining: qualitative information on the main findings reported by studies, quantitative data on the performance of these studies, detailed quantitative analysis of the 206 models (or model variants) reported in 19 studies which report (or we can calculate from what is reported) precision, recall, and f-measure performance data.

This paper is organized as follows: In the next section, we present our systematic literature review methodology. In Section 3, we present our criteria developed to assess whether or not a study reports sufficient contextual and methodological detail to enable us to synthesize a particular study. Section 4 shows the results of applying our assessment criteria to 208 studies. Section 5 reports the results of extracting data from the 36 studies which satisfy

our assessment criteria. Section 6 synthesizes our results and Section 7 discusses the methodological issues associated with fault prediction studies. Section 8 identifies the threats to validity of this study. Finally, in Section 9 we summarize and present our conclusions.

2 METHODOLOGY

We take a systematic approach to reviewing the literature on the prediction of faults in code. Systematic literature reviews are well established in medical research and increasingly in software engineering. We follow the systematic literature review approach identified by Kitchenham and Charters [3].

2.1 Research Questions

The aim of this systematic literature review is to analyze the models used to predict faults in source code. Our analysis is based on the research questions in Table 1.

2.2 Inclusion Criteria

To be included in this review, a study must be reported in a paper published in English as either a journal paper or conference proceedings. The criteria for studies to be included in our SLR are based on the inclusion and exclusion criteria presented in Table 2.

Before accepting a paper into the review, we excluded repeated studies. If the same study appeared in several publications, we included only the most comprehensive or most recent.

2.3 Identification of Papers

Included papers were published between January 2000 and December 2010. Our searches for papers were completed at the end of May 2011 and it is therefore unlikely that we missed any papers published in our time period as a result of publication time lags. There were four elements to our searches:

1. Key word searching using the search engines: ACM Digital Library, IEEEExplore, and the ISI Web of Science. These search engines covered the vast

TABLE 2
Inclusion and Exclusion Criteria

Inclusion criteria (a paper must be...)	Exclusion criteria (a paper must not be...)
<ul style="list-style-type: none"> - An empirical study - Focused on predicting faults in units of a software system - Faults in code is the main output (dependent variable) 	<ul style="list-style-type: none"> - Focused on: testing, fault injection, inspections, reliability modelling, aspects, effort estimation, debugging, faults relating to memory leakage, nano-computing, fault tolerance. - About the detection or localisation of existing individual known faults.

TABLE 3
Paper Selection and Validation Process

Selection Process	# of papers	Validation
Papers extracted from databases, conferences and journals	2,073	80 random papers independently classified by 3 researchers
Sift based on title and abstract	-1,762 rejected	Fair/good inter-rater agreement on first sift (k statistic test)
Full papers considered for review	311 primary 80 secondary	Each paper is read in full and 80 secondary papers are identified from references
Rejected on full reading	-185 rejected	Papers are rejected on the basis that they do not answer our research questions
Comparison to Catal and Diri's review	2 (our searches missed)	
Papers accepted for the review	208 papers	

majority of software engineering publications and the search string we used is given in Appendix A.

2. An issue-by-issue manual reading of paper titles in relevant journals and conferences. The journals and conferences searched are shown in Appendix B. These were chosen as highly relevant software engineering publications found previously to be good sources of software engineering research [4].
3. A manual search for publications from key authors using DBLP.³ These authors were selected as appearing most frequently in our list of papers: Khoshgof-taar, Menzies, Nagappan, Ostrand, and Weyuker.
4. The identification of papers using references from included studies.

Table 3 shows that our initial searches elicited 2,073 papers. The title and abstract of each was evaluated and 1,762 were rejected as not relevant to fault prediction. This process was validated using a randomly selected 80 papers from the initial set of 2,073. Three researchers separately interpreted and applied the inclusion and exclusion criteria to the 80 papers. Pairwise interrater reliability was measured across the three sets of decisions to get a fair/good agreement on the first iteration of this process. Based on the disagreements, we clarified our inclusion and exclusion criteria. A second iteration resulted in 100 percent agreement between the three researchers.

We read the remaining 311 papers in full. This resulted in a further 178 papers being rejected. An additional 80 secondary papers were identified from references and, after being read in full, accepted into the included set. We also included two extra papers from Catal and Diri's [2] review which overlapped our timeframe. Our initial searches omitted these two of Catal and Diri's papers as their search terms included the word "quality." We did not include this word in our searches as it generates a very high false positive rate. This process resulted in the 208 papers included in this review.

3 ASSESSING THE SUITABILITY OF PAPERS FOR SYNTHESIS

The previous section explained how we included papers which both answered our research questions and satisfied our inclusion criteria. This section describes how we identified a subset of those papers as suitable from which to extract data and synthesize an overall picture of fault prediction in software engineering. We then describe the extraction and synthesis process.

3.1 The Assessment Criteria

Our approach to identifying papers suitable for synthesis is motivated by Kitchenham and Charter's [3] notion of a quality check. Our assessment is focused specifically on identifying only papers reporting sufficient information to allow synthesis across studies in terms of answering our research questions. To allow this, a basic set of information must be reported in papers. Without this it is difficult to properly understand what has been done in a study and equally difficult to adequately contextualize the findings reported by a study. We have developed and applied a set of criteria focused on ensuring sufficient contextual and methodological information is reported in fault prediction studies. Our criteria are organized into four phases described below.

Phase 1: Establishing that the study is a prediction study.

In this SLR it is important that we consider only models which actually do some form of prediction. Some studies which seem to be reporting prediction models actually turn out to be doing very little prediction. Many of these types of studies report correlations between metrics and faults. Such studies only indicate the propensity for building a prediction model. Furthermore, a model is only doing any prediction if it is tested on unseen data (i.e., data that were not used during the training process) [S112]. To be considered a prediction model it must be trained and tested on different data [6]. Table 4 shows the criteria we apply to assess whether a study is actually a prediction study.

3. <http://www.informatik.uni-trier.de/~ley/db/>.

TABLE 4
Prediction Criteria

Prediction criteria	Criteria definitions	Why the criteria is important
Is a prediction model reported?	The study must report some form of prediction. Not just report a correlation study of the relationship between faults and independent variables.	Such studies provide useful insights into observed patterns of faults but do not present prediction models as such. Nor do they validate their findings using unseen data. We do not therefore take these studies forward for synthesis.
Is the prediction model tested on unseen data?	The prediction model must be developed and tested on different data. This means that some form of hold-out or cross validation is necessary.	The performance of predictions based only on training data gives us no information on which to judge the performance of how such models generalise to new data. Such studies are therefore not taken forward for synthesis.

Table 4 shows that a study can pass this criterion as long as they have separated their training and testing data. There are many ways in which this separation can be done. Holdout is probably the simplest approach, where the original dataset is split into two groups comprising: {training set, test set}. The model is developed using the training set and its performance is then assessed on the test set. The weakness of this approach is that results can be biased because of the way the data have been split. A safer approach is often n -fold cross validation, where the data are split into n groups $\{g_1 \dots g_n\}$. Ten-fold cross validation is very common, where the data are randomly split into 10 groups, and 10 experiments carried out. For each of these experiments, one of the groups is used as the testing set, and all others combined are used as the training set. Performance is then typically reported as an average across all 10 experiments. M-N fold cross validation adds another step by generating M different N-fold cross validations, which increases the reliability of the results and reduces problems due to the order of items in the training set.

Stratified cross validation is an improvement to this process, and keeps the distribution of faulty and nonfaulty data points approximately equal to the overall class distribution in each of the n bins. Although there are stronger and weaker techniques available to separate training and testing data, we have not made a judgment on this and have accepted any form of separation in this phase of assessment.

Phase 2: Ensuring sufficient contextual information is reported.

We check that basic contextual information is presented by studies to enable appropriate interpretation of findings. A lack of contextual data limits the user's ability to: interpret a model's performance, apply the model appropriately, or repeat the study. For example, a model may have been built using legacy systems with many releases over a long time period and has been demonstrated to perform well on these systems. It may not then make sense to rely on this model for a new system where the code has only recently been developed. This is because the number and type of faults in a system are thought to change as a system evolves [S83]. If the maturity of the system on which the model was built is not reported, this severely limits a model user's ability to understand the conditions in which the model performed well and to select this model specifically for legacy systems. In this situation the model could be applied to newly developed systems with disappointing predictive performance.

The contextual criteria we applied are shown in Table 5 and are adapted from the context checklist developed by Petersen and Wohlin [7]. Our context checklist also overlaps with the 40 project characteristics proposed by Zimmermann et al. [S208] as being relevant to understanding a project sufficiently for cross project model building (it was impractical for us to implement all 40 characteristics as none of our included studies report all 40).

Context data are particularly important in this SLR as it is used to answer Research Question 1 and interpret our overall findings on model performance. We only synthesize papers that report all the required context information as listed in Table 5. Note that studies reporting several models based on different datasets can pass the criteria in this phase if sufficient contextual data are reported for one or more of these models. In this case, data will only be extracted from the paper based on the properly contextualized model.

Phase 3: Establishing that sufficient model building information is reported.

For a study to be able to help us to answer our research questions it must report its basic model building elements. Without clear information about the independent and dependent variables used as well as the modeling technique, we cannot extract sufficient data to allow synthesis. Table 6 describes the criteria we apply.

Phase 4: Checking the model building data.

Data used are fundamental to the reliability of models. Table 7 presents the criteria we apply to ensure that studies report basic information on the data they used.

In addition to the criteria we applied in Phases 1 to 4, we also developed more stringent criteria that we did not apply. These additional criteria relate to the quality of the data used and the way in which predictive performance is measured. Although we initially intended to apply these, this was not tenable because the area is not sufficiently mature. Applying these criteria would have resulted in only a handful of studies being synthesized. We include these criteria in Appendix C as they identify further important criteria that future researchers should consider when building models.

3.2 Applying the Assessment Criteria

Our criteria have been applied to our included set of 208 fault prediction studies. This identified a subset of 36 finally included studies from which we extracted data and on which our synthesis is based. The initial set of 208 included papers was divided between the five authors. Each paper was assessed by two authors independently

TABLE 5
Context Criteria

Contextual criteria	Criteria definitions	Why the criteria is important
Source of data	The source of the system data on which the study is based must be given. For example whether the system data is industrial, open source, NASA, Promise. If NASA/Promise data is used the names of the datasets used must be given. Studies using data that is in the public domain, the context of which is accessible via the public domain, need not explicitly report all these criteria to pass this phase. However the version studied must be specified to enable access to contextual data.	Different models may perform differently when applied to different data sets; for example some models may perform better on OS data than industrial data. It is therefore essential for synthesis that we can establish the source of the data.
Maturity	Some indication of the maturity of the system being studied must be reported. Readers must be able to generally determine whether the system is a mature system with many releases which has been in the field for many years, or whether it is a relatively newly developed system, or whether the system has yet to be released.	The age of a system has a significant impact on how it behaves. Especially in terms of the faults in the system. Many factors contribute to why the age of the system impacts on faults in the system, including the amount of change the system has undergone. This means that some models are likely to perform better than others on newly developed as opposed to legacy systems. It is therefore essential that we can establish the maturity of the system(s) on which the model was based, as without this it is difficult to correctly interpret study findings for synthesis.
Size in KLOC	An indication of the size of the system being studied must be given in KLOC. The overall size of the system will suffice, i.e. it is not necessary to give individual sizes of each component of the system being studied, even if only sub-sets of the system are used during prediction. Size indicated by measures other than KLOC are not acceptable (e.g. number of classes) as there are great variations in the KLOC of such other measures.	The size of systems is likely to impact on the behaviour of systems. Consequently, the faults in a system may be different in small as opposed to large systems. This means that it is also likely that different types of models will perform differently on systems of different sizes. It is therefore essential that we can establish the size of the system(s) on which the model was built as without this it is difficult to correctly interpret study findings for synthesis. Using KLOC does have limitations. The definitions of KLOC can vary and KLOC can be language dependent [8].
Application domain	A general indication of the application domain of the system being studied must be given, e.g. telecoms, customer support, etc.	Some models are likely to be domain specific. Different domains apply different development practices and result in different faults. It is therefore important that domain information is given so that this factor can be taken into account when model performance is evaluated. It is therefore essential that we establish the domain of the system(s) on which the model was built, as without this it is difficult to correctly interpret study findings for synthesis.
Programming language	The programming language(s) of the system being studied must be given.	Different languages may result in different faults. In particular it may be that OO languages perform differently from procedural languages. This makes it likely that some models will perform better for some languages. It is therefore essential that we can establish the language of the system(s) on which the model was built as without this it is difficult to correctly interpret study findings for synthesis.

TABLE 6
Model Building Criteria

Model building criteria	Criteria definitions	Why the criteria is important
Are the independent variables clearly reported?	The basis on which the model is making predictions must be clear and explicit. For example independent variables (or predictor variables) could include: static code metrics (complexity etc.), code churn metrics, previous fault metrics, etc.	In any experimental work which shows the performance of a method it is essential that the independent variables tested are explicitly identified. Without this, confidence in the work is significantly lowered and it is difficult to evaluate the impact of those variables across studies during synthesis.
Is the dependent variable clearly reported?	It must be distinguishable whether studies are predicting faults in terms of whether a module is fault prone or not fault prone (i.e. using a categorical dependent variable) or in terms of the number of faults in a code unit (i.e. using a continuous dependent variable). Some continuous studies additionally report ranked results (i.e. the identification of the faultiest 20% of files).	In any experimental work it is essential that the dependent variables are explicitly identified. Without this, confidence in the work is significantly lowered. Furthermore, the evaluation of fault prediction models is related to whether the dependent variable is categorical or continuous. It is therefore essential for synthesis that this information is clear.
Is the granularity of the dependent variable reported?	The unit of code granularity of predictions must be reported. For example fault predictions in terms of faults per module, per file, per package etc. These terms may be used differently by different authors, e.g. 'module' is often used to mean different code units by different authors. Studies must indicate the code unit being used, i.e. if the term 'module' is used, readers must be able to work out what unit of code is being defined as a module.	It is difficult to directly compare the performance of one model reporting faults per file to another model reporting faults per method. Furthermore it may be that studies reporting faults at a file level are more able to perform well than studies reporting at a method level. The granularity of the dependent variable must therefore be taken into account during synthesis and so an indication of the unit of fault granularity must be reported.
Is the modelling technique used reported?	It must be clear what modelling technique is being used, e.g. linear regression, decision trees, etc. It is not acceptable to present results from a tool based on a model that is not discussed in the paper.	Different modelling techniques are likely to perform differently in different circumstances. A study must report the modelling technique used, as we cannot examine the impact of method on performance without this information.

TABLE 7
Data Criteria

Data criteria	Criteria definitions	Why the criteria is important
Is the fault data acquisition process described?	The process by which the fault data was obtained must be described. A high level indication of this will suffice, e.g. obtained from the CVS records. However at least an overview of how the data was obtained must be provided. If the data has been obtained from a third party (e.g. NASA), some indication of how that data was obtained by that third party must be given or referenced or available in the public domain.	In order to have confidence in the data on which a model was built it is necessary to have some information on how the data was collected. Data is fundamental to the quality of the models built and the collection process has a huge impact on the quality of that data. It is not possible to have any confidence in a study where the data seems to have come from thin air. Collecting software engineering data of any sort is difficult and error-prone. We do not synthesise papers which give no indication of how the data was collected. In an ideal world studies would also indicate the development point at which fault data collection starts and stops. This is because fault data may be unreliable for several reasons: fault data may be extracted from a different version of the system to which the independent variable relates; faults are likely to emerge after fault counting ends, and data collected at the start of system development is unlikely to be stable or complete.
Is the independent variable data acquisition process described?	Some indication of how the independent variable data (e.g. static code data) was obtained should be given. For example the static analysis tools used should be reported or the process by which code churn data collected should be described. This does not need to be at a great level of detail, just an indication given. If the data has been obtained from a third party (e.g. NASA), some indication of how that data was obtained by that third party must be given, referenced or available in the public domain.	As above.
For categorical studies, has the number of faulty versus non-faulty units on which the model has been trained and tested on been reported?	The balance of faulty versus non-faulty units used for training and testing must be reported. For studies using open source systems it is not essential to report this (though preferable) as this data should be possible to identify from the public data source.	The balance of faulty versus non-faulty units used for training and testing can affect the reliability of some performance measures (see Appendix F). It is essential that class distributions are reported (i.e. number of faulty and number of non-faulty units in the data used). This makes it possible to appropriately interpret performances reported using these measures. We use this information for our synthesis of categorical studies. As where precision and recall are not reported by such studies we re-compute an approximation of it. The faulty/non-faulty balance of data is often needed in this calculation.

(with each author being paired with at least three other authors). Each author applied the assessment criteria to between 70 and 80 papers. Any disagreements on the assessment outcome of a paper were discussed between the two authors and, where possible, agreement established between them. Agreement could not be reached by the two authors in 15 cases. These papers were then given to another member of the author team for moderation. The moderator made a final decision on the assessment outcome of that paper.

We applied our four phase assessment to all 208 included studies. The phases are applied sequentially. If a study does not satisfy all of the criteria in a phase, then the evaluation is stopped and no subsequent phases are applied to the study. This is to improve the efficiency of the process as there is no point in assessing subsequent criteria if the study has already failed the assessment. This does have the limitation that we did not collect information on how a paper performed in relation to all assessment criteria. So if a paper fails Phase 1, we have no information on how that paper would have performed in Phase 4.

This assessment process was piloted four times. Each pilot involved three of the authors applying the assessment to 10 included papers. The assessment process was refined as a result of each pilot.

We developed our own MySQL database system to manage this SLR. The system recorded full reference details

and references to pdfs for all papers we identified as needing to be read in full. The system maintained the status of those papers as well as providing an online process to support our assessments of 208 papers. The system collected data from all authors performing assessments. It also provided a moderation process to facilitate identifying and resolving disagreements between pairs of assessors. The system eased the administration of the assessment process and the analysis of assessment outcomes. All data that were extracted from the 36 papers which passed the assessment is also recorded on our system. An overview of the system is available from [9] and full details are available from the third author.

3.3 Extracting Data from Papers

Data addressing our three research questions was extracted from each of the 36 finally included studies which passed all assessment criteria. Our aim was to gather data that would allow us to analyze predictive performance within individual studies and across all studies. To facilitate this, three sets of data were extracted from each study:

1. **Context data.** Data showing the context of each study were extracted by one of the authors. This data give the context in terms of: the source of data studied and the maturity, size, application area, and programming language of the system(s) studied.

2. **Qualitative data.** Data related to our research questions were extracted from the findings and conclusions of each study. This was in terms of what the papers reported rather than on our own interpretation of their study. This data supplemented our quantitative data to generate a rich picture of results within individual studies.

Two authors extracted qualitative data from all 36 studies. Each author extracted data independently and compared their findings to those of the other author. Disagreements and omissions were discussed within the pair and a final set of data agreed upon.

3. **Quantitative data.** Predictive performance data were extracted for every individual model (or model variant) reported in a study. The performance data we extracted varied according to whether the study reported their results via categorical or continuous dependent variables. Some studies reported both categorical and continuous results. We extracted only one of these sets of results, depending on the way in which the majority of results were presented by those studies. The following is an overview of how we extracted data from categorical and continuous studies.

Categorical studies. There are 23 studies reporting categorical dependent variables. Categorical studies report their results in terms of predicting whether a code unit is likely to be fault prone or not fault prone. Where possible we report the predictive performance of these studies using precision, recall, and f-measure (as many studies report both precision and recall, from which an f-measure can be calculated). F-measure is commonly defined as the harmonic mean of precision and recall, and generally gives a good overall picture of predictive performance.⁴ We used these three measures to compare results across studies and, where necessary, we calculate and derive these measures from those reported (Appendix E explains how we did this conversion and shows how we calculated f-measure). Standardizing on the performance measures reported allows comparison of predictive performances across studies. Lessmann et al. [S97] recommend the use of consistent performance measures for cross-study comparison; in particular, they recommend use of Area Under the Curve (AUC). We also extract AUC where studies report this. Appendix D summarizes the measurement of predictive performance.

We present the performance of categorical models in boxplots. Box plots are useful for graphically showing the differences between populations. They are useful for our results as they make no assumptions about the distribution of the data presented. These boxplots present the precision, recall, and f-measure of studies according to a range of model factors. These factors are related to the research questions presented at the beginning of Section 2; an example is a boxplot showing model performance relative to the modeling technique used.

4. Menzies et al. [10] claim that values of precision vary greatly when used with models applied to different datasets. However, reporting precision and recall via an f-measure effectively evaluates classifier performance, even in highly imbalanced domains [11], [12].

Continuous studies. There are 13 studies reporting continuous dependent variables. These studies report their results in terms of the number of faults predicted in a unit of code. It was not possible to convert the data presented in these studies into a common comparative measure; we report the individual measures that they use. Most measures reported by continuous studies are based on reporting an error measure (e.g., Mean Standard Error (MSE)), or measures of difference between expected and observed results (e.g., Chi Square). Some continuous studies report their results in ranking form (e.g., top 20 percent of faulty units). We extract the performance of models using whatever measure each study used.

Two authors extracted quantitative data from all 36 studies. A pair approach was taken to extracting this data since it was a complex and detailed task. This meant that the pair of authors sat together identifying and extracting data from the same paper simultaneously.

3.4 Synthesizing Data across Studies

Synthesizing findings across studies is notoriously difficult and many software engineering SLRs have been shown to present no synthesis [13]. In this paper, we have also found synthesizing across a set of disparate studies very challenging. We extracted both quantitative and qualitative data from studies. We intended to meta-analyze our quantitative data across studies by combining precision and recall performance data. However, the studies are highly disparate in terms of both context and models. Meta-analyzing this quantitative data may generate unsafe results. Such a meta-analysis would suffer from many of the limitations in SLRs published in other disciplines [14].

We combined our qualitative and quantitative data to generate a rich picture of fault prediction. We did this by organizing our data into themes based around our three research questions (i.e., context, independent variables, and modeling techniques). We then combined the data on each theme to answer our research questions. This synthesis is presented in Section 6.

4 RESULTS OF OUR ASSESSMENT

This section presents the results from applying our assessment criteria (detailed in Tables 4, 5, 6, and 7) to establish whether or not a paper reports sufficient contextual and methodological detail to be synthesized. The assessment outcome for each study is shown at the end of its reference in the list of included studies.

Table 8 shows that only 36 of our initially included 208 studies passed all assessment criteria.⁵ Of these 36 finally included studies, three are relatively short [S116], [S110], and [S164]. This means that it is possible to report necessary contextual and methodological detail concisely without a significant overhead in paper length. Table 8 also shows that 41 papers failed at phase 1 of the assessment because they did not report prediction models as such. This includes studies that only present correlation studies or

5. These papers are [S8], [S9], [S10], [S11], [S18], [S21], [S29], [S31], [S32], [S37], [S51], [S56], [S69], [S73], [S74], [S76], [S83], [S86], [S92], [S98], [S109], [S110], [S116], [S117], [S118], [S120], [S122], [S127], [S133], [S135], [S154], [S160], [S163], [S164], [S190], [S203].

TABLE 8
Results of Applying Assessment Criteria

Number of papers passed	Number of papers failed					
	Phase 1: Prediction	Phase 2: Context	Phase 3: Model	Phase 4: Data	Other reasons	Total failed
36	41	114	2	13	2	173

models that were not tested on data unseen during training. This is an important finding as it suggests that a relatively high number of papers reporting fault prediction are not really doing any prediction (this finding is also reported by [6]).

Table 8 also shows that 13 studies provided insufficient information about their data. Without this it is difficult to establish the reliability of the data on which the model is based. Table 8 also shows that a very high number of studies (114) reported insufficient information on the context of their study. This makes it difficult to interpret the results reported in these studies and to select an appropriate model for a particular context. Several studies passing all of our criteria anonymized their contextual data, for example, [S109] and [S110]. Although these studies gave full contextual details of the systems they used, the results associated with each were anonymized. This meant that it was impossible to relate specific fault information to specific systems. While a degree of commercial confidentiality was maintained, this limited our ability to analyze the performance of these models.

Of the 114 studies which did not report sufficient context information, 58 were based on NASA data (located in NASA MDP or PROMISE). This is because we could find no information about the maturity of the systems on which the NASA data are based. Maturity information is not given in either the MDP or PROMISE repository documentation and no included paper provided any maturity information. Turham et al. [15] report that the NASA data are from numerous NASA contractors for an array of projects with a wide range of reuse. This suggests that a range of maturities might also be represented in these datasets. No clear insight is given into whether particular datasets are based on systems developed from untested, newly released, or legacy code based on many releases. The only three studies using NASA data which passed the context phase of the assessment were those which also used other datasets for which full context data are available (the NASA-based models were not extracted from these studies). Whether a study uses NASA data (sourced from MDP or PROMISE) is shown at the end of its reference in the list of included studies.

Table 8 also shows that two studies failed the assessment due to the “other” reasons reported in Table 9.

5 RESULTS EXTRACTED FROM PAPERS

This section presents the results we extracted from the 36 papers that passed all of our assessment criteria. The full set of data extracted from those papers are contained in our online Appendix (<https://bugcatcher.stca.herts.ac.uk/slr2011/>). This online Appendix consists of the following.

1. **Context of study table.** For each of the 36 studies, the context of the study is given in terms of: the aim of the study together with details of the system(s) used in the study (the application area(s), the system(s), maturity, and size(s)).
2. **Categorical models table.** For each study reporting categorical results, each model is described in terms of the: independent variable(s), the granularity of the dependent variable, the modeling technique(s), and the dataset(s) used. This table also reports the performances of each model using precision, recall, f-measure, and (where given by studies) AUC. Some studies present many models or model variants, all of which are reported in this table.
3. **Continuous models table.** For each study reporting continuous results (including those reporting ranking results) the same information describing their model(s) is presented as for categorical models. However, the performance of each continuous model is reported in terms of either: the error measure, the measure of variance, or the ranked results (as reported by a study).
4. **Qualitative data table.** For each study a short summary of the main findings reported by authors is presented.

The remainder of this section contains boxplots illustrating the performance of the models in relation to various model factors (e.g., modeling technique used, independent variable used, etc.). These factors are related to the research questions that we posed at the beginning of Section 2. The boxplots in this section set performance against individual model factors (e.g., modeling technique used). This is a simplistic analysis, as a range of interacting factors are likely to underpin the performance of a model. However, our results indicate areas of promising future research.

TABLE 9
Issues with the Measurement of Performance

Paper number	Performance measurement issues
[[141]]	The model was optimised on the test set possibly resulting in inflated performance as shown by the 100% accuracy.
[[52]]	The paper does not report the error rate of any of the models it presents. Only the difference between the error rates reported by each model reported is given. However the error rate could be very high for all models, but the difference between each model could be very small giving the impression that the models are all working well.

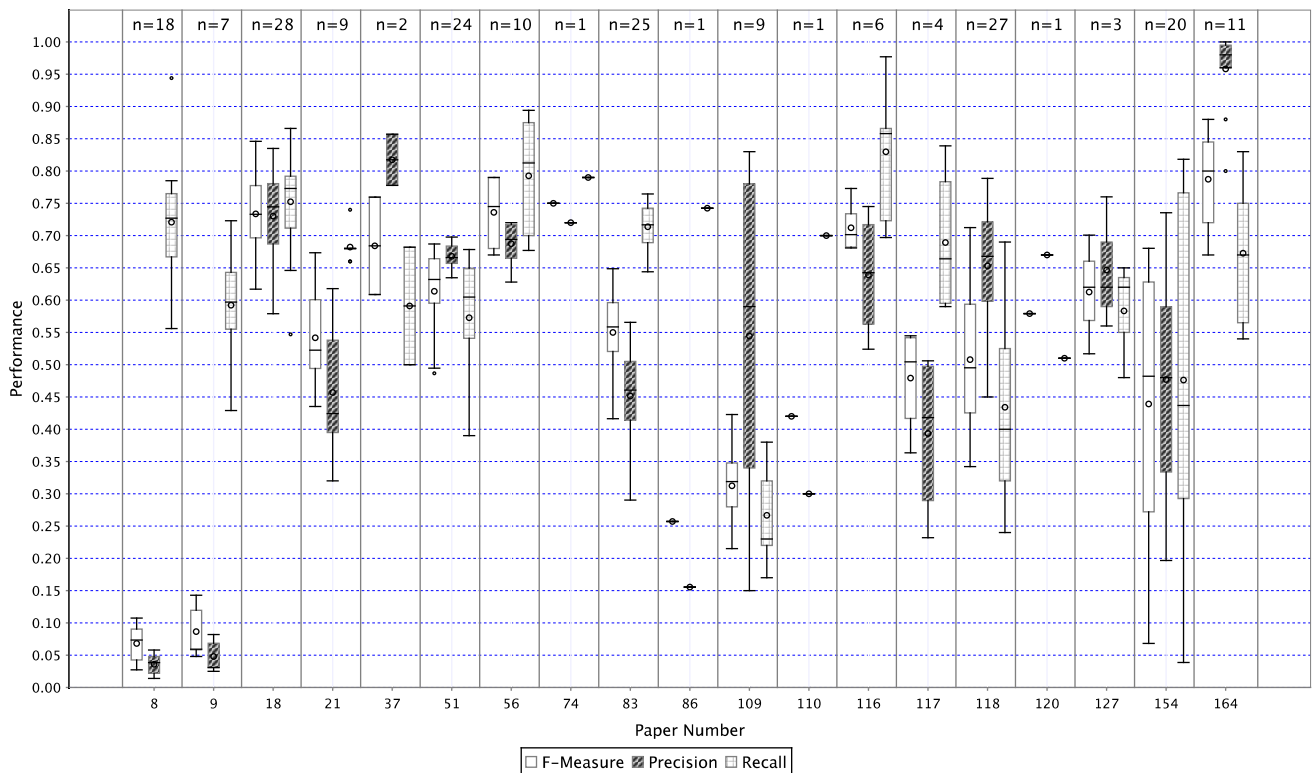


Fig. 1. Performances of the models reported in each of the categorical studies.

The boxplots represent models reporting only categorical results for which precision, recall, and f-measure were either reported or could be calculated by us. Such models are reported in 19 of the 23 categorical studies (of the remaining four, three report AUC). We are unable to present boxplots for the 13 studies using continuous data as the measures used are not comparable or convertible to comparable measures.

Each boxplot includes data only where at least three models have used a particular factor (e.g., a particular independent variable like LOC). This means that the numbers (n) at the top of the boxplots will not add up to the same number on every plot, as factors used in less than three studies will not appear; the total of n s will therefore vary from one boxplot to the next. The boxplots contain performance data based on precision, recall, and f-measure. This is for all categorical models and model variants presented by each study (206 models or model variants). Some studies present many model variants while others present only one model. We also created boxplots of only the best results from each study. These boxplots did not change the pattern of good performances but only presented limited information about poor performances. For that reason, we do not include these “best only” boxplots.

5.1 Performances of Models Reported in Individual Studies

Fig. 1 is a boxplot of the performances of all the models reported by each of the 19 categorical papers (full details of which can be found in the online Appendix). For each individual paper, f-measure, precision, and recall is reported. Fig. 1 shows that studies report on many models or variants of models, some with a wide range of performances

(the details of these can be found in the Models Table in the online Appendix (<https://bugcatcher.stca.herts.ac.uk/slr2011/>)). For example, Schröter et al. [S154] present 20 model variants with a wide range of precision, recall, and f-measure. Many of these variants are not particularly competitive; the most competitive models that Schröter et al. [S154] report are based on training the model on only the faultiest parts of the system. This is a promising training technique and a similar technique has also been reported to be successful by Zhang et al. [S200]. Bird et al. [S18] report 28 model variants with a much smaller range of performances, all of which are fairly competitive. Fig. 1 also shows the performance tradeoffs in terms of precision and recall made by some models. For example, Bird et al. [S18] report consistent precision and recall, whereas Moser et al. [S118] and Shivaji et al. [S164] report performances where precision is much higher than recall.

Fig. 1 also shows that some models seem to be performing better than others. The models reported by Shivaji et al. [S164], based on Naive Bayes, performed extremely competitively. In general Naive Bayes performed relatively well, see Fig. 8. However, Shivaji et al. [S164] also used a good modeling process, including feature selection and appropriate measures derived during model training. In addition, their dataset contained a relatively large proportion of faulty components, making it fairly balanced. This may improve performance by providing many examples of faults from which the modeling technique can train. There are many good aspects of this study that mean it is likely to produce models which perform well.

On the other hand, the performance of Arisholm et al.’s models [S8], [S9] are low in terms of precision but

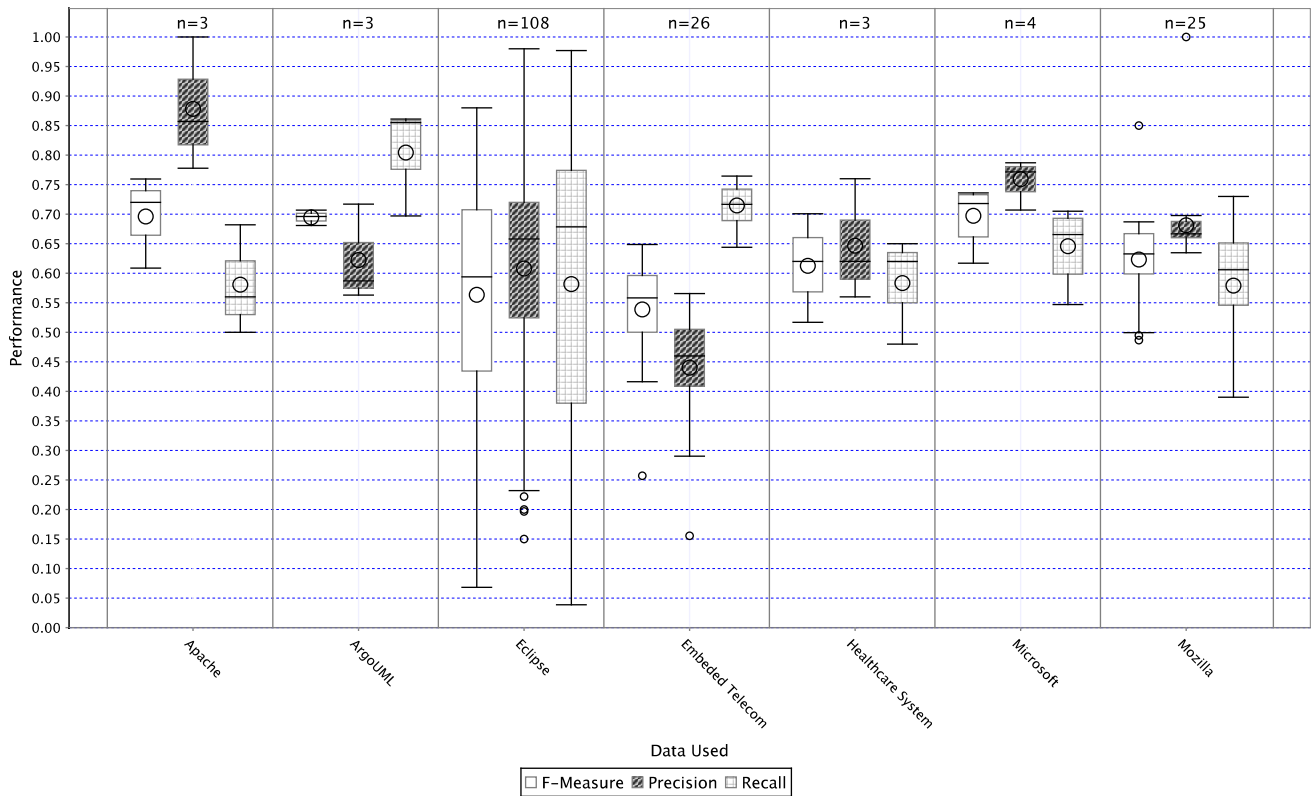


Fig. 2. Data used in models.

competitive in terms of recall. The two Arisholm et al. studies are different but use the same datasets. This low precision is reportedly because of the sampling method used to address the imbalance of the data used. Though the datasets used are also small relative to those used in other studies (148 KLOC), Arisholm et al.'s studies [S8], [S9] are interesting as they also report many good modeling practices and in some ways are exemplary studies. But they demonstrate how the data used can impact significantly on the performance of a model. It is also essential that both high and low performances be reported, as it is only by identifying these that our overall understanding of fault prediction will improve. The boxplots in the rest of this section explore in more detail aspects of models that may underpin these performance variations. Because the performances of Arisholm et al.'s models [S8], [S9] are very different from those of the other studies, we have removed them from the rest of the boxplots. We have treated them as outliers which would skew the results we report in other boxplots.

5.2 Performances in Relation to Context Factors

Fig. 2 shows the datasets used in the studies. It shows that 108 models reported in the studies are based on data from Eclipse. Eclipse is very well studied, probably because the fault data are easy to access and its utility has been well proven in previous studies. In addition, data already extracted from Eclipse are available from Saarland University (<http://www.st.cs.uni-saarland.de/softevo/bug-data/eclipse/>) and PROMISE (<http://promisedata.org/>). Fig. 2 shows that there is a wide variation in model performance using Eclipse. Fig. 2 also suggests that it may be more difficult to build models for some systems than for

others. For example, the models built for embedded telecoms systems are not particularly competitive. This may be because such systems have a different profile of faults with fewer postdelivery faults relative to other systems. Developers of such systems normally prioritize reducing postdelivery faults as their embedded context makes fixing them comparatively expensive [S83].

Fig. 3 shows how models have performed relative to the size of systems on which they are based. Eclipse is the most common system used by studies. Consequently, Fig. 3 shows only the size of versions of Eclipse in relation to model performance. Fig. 3 suggests that as the size of a system increases, model performance seems to improve. This makes sense as models are likely to perform better given more data.

Fig. 4 shows the maturity of systems used by studies relative to the performance of models. The Context Table in the online Appendix shows how systems have been categorized in terms of their maturity. Fig. 4 shows that no immature systems are used by more than two models in this set of studies (i.e., where $n \geq 3$).⁶ There seems to be little difference between the performance of models using mature or very mature systems. This suggests that the maturity of systems may not matter to predictive performance.⁷ This finding may be linked to the finding we report on size. It may be that what was previously believed about the

6. An exception to this is found in studies [S11], [S133], where immature systems are used with promising performances reported (see the online Appendix for full details).

7. This may mean that it is not important to report maturity when studies describe their context (many more studies would have passed our assessment had that been the case). However, much more data on maturity is needed before firm conclusions can be drawn.

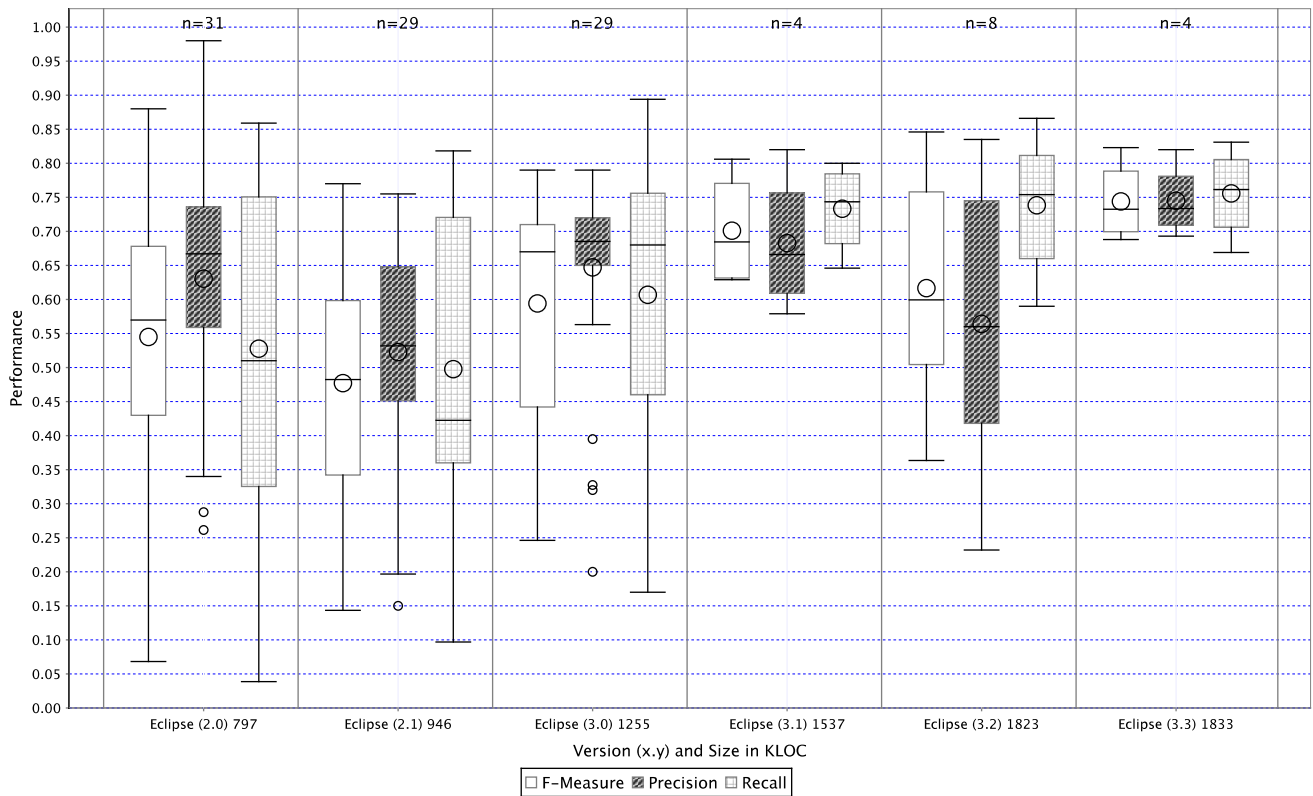


Fig. 3. The size of the datasets used for Eclipse.

importance of maturity was actually about size, i.e., maturity is a surrogate for size. Indeed, there is a significant relationship between size and maturity in the data we report here. However, we do not have enough data to draw firm conclusions as the data we analyze contain no studies using immature systems. More research is needed to test for possible association between maturity and size and whether data extracted from immature systems can be used as a basis for reliable fault prediction.

Fig. 5 shows the language used in the systems studied in relation to the performance of models. We present only studies reporting the use of either Java or C/C++. There are

several single studies using other languages which we do not report. Fig. 5 suggests that model performance is not related to the language used.

Fig. 6 shows model performance relative to the granularity of dependent variables (e.g., whether fault prediction is at the class or file level). It shows no clear relationship between granularity and performance. It does not seem to be the case that higher granularity is clearly related to improved performance. Models reporting at “other” levels of granularity seem to be performing most consistently. These tend to be high levels of granularity defined specifically by individual studies (e.g., Nagappan et al. [S120]).

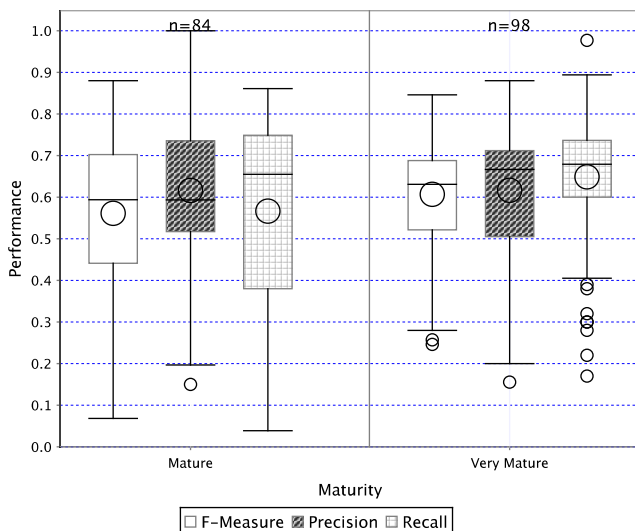


Fig. 4. The maturity of the systems used.

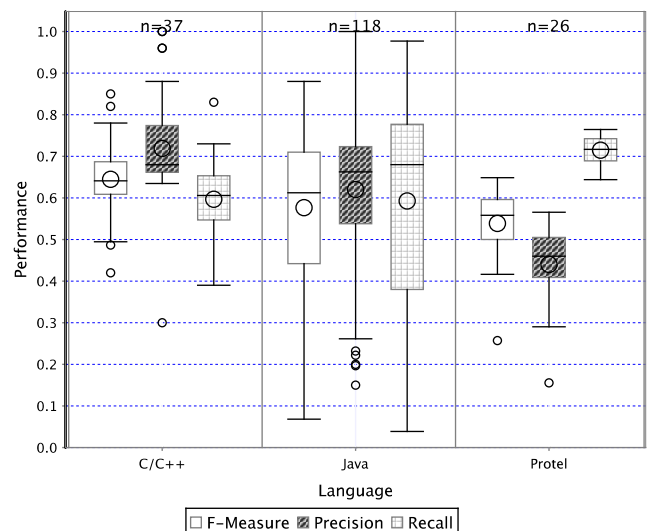
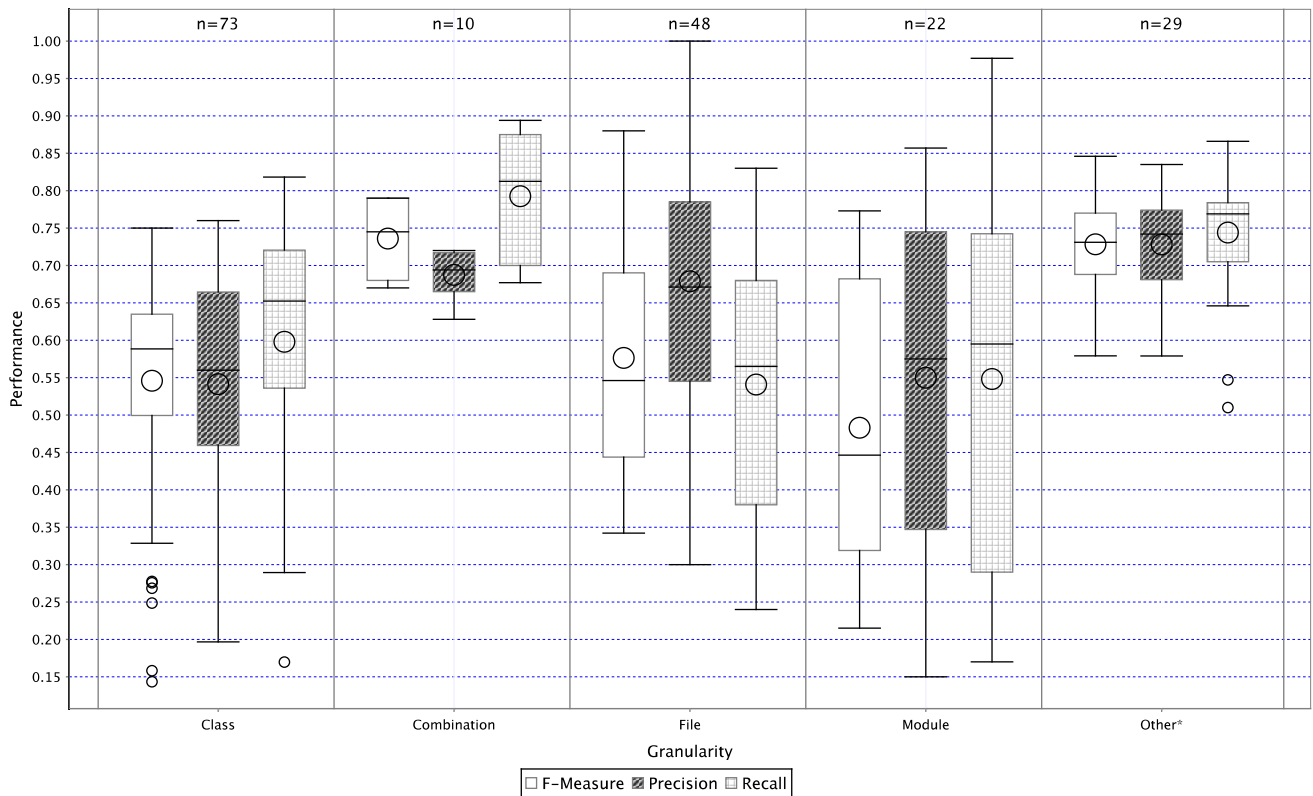


Fig. 5. The language used.



*For example plug-ins, binaries

Fig. 6. The granularity of the results.

5.3 Performance in Relation to Independent Variables

Fig. 7 shows model performance in relation to the independent variables used. The Categorical Models Table in the online Appendix shows how independent variables as expressed by individual studies have been categorized in relation to the labels used in Fig. 7. It shows that there is variation in performance between models using different independent variables. Models using a wide combination of metrics seem to be performing well. For example, models using a combination of static code metrics (scm), process metrics, and source code text seem to be performing best overall (e.g., Shivaji et al. [S164]). Similarly Bird et al.'s study [S18], which uses a wide combination of socio-technical metrics (code dependency data together with change data and developer data), also performs well (though the results from Bird et al.'s study [S18] are reported at a high level of granularity). Process metrics (i.e., metrics based on changes logged in repositories) have not performed as well as expected. OO metrics seem to have been used in studies which perform better than studies based only on other static code metrics (e.g., complexity-based metrics). Models using only LOC data seem to have performed competitively compared to models using other independent variables. Indeed, of these models using only metrics based on static features of the code (OO or SCM), LOC seems as good as any other metric to use. The use of source code text seems related to good performance. Mizuno et al.'s studies [S116], [S117] have used only source

code text within a novel spam filtering approach to relatively good effect.

5.4 Performance in Relation to Modeling Technique

Fig. 8 shows model performance in relation to the modeling techniques used. Models based on Naive Bayes seem to be performing well overall. Naive Bayes is a well understood technique that is in common use. Similarly, models using Logistic Regression also seem to be performing well. Models using Linear Regression perform not so well, though this technique assumes that there is a linear relationship between the variables. Studies using Random Forests have not performed as well as might be expected (many studies using NASA data use Random Forests and report good performances [S97]). Fig. 8 also shows that SVM (Support Vector Machine) techniques do not seem to be related to models performing well. Furthermore, there is a wide range of low performances using SVMs. This may be because SVMs are difficult to tune and the default Weka settings are not optimal. The performance of models using the C4.5 technique is fairly average. However, Arisholm et al.'s models [S8], [S9] used the C4.5 technique (as previously explained, these are not shown as their relatively poor results skew the data presented). C4.5 is thought to struggle with imbalanced data [16] and [17] and this may explain the performance of Arisholm et al.'s models.

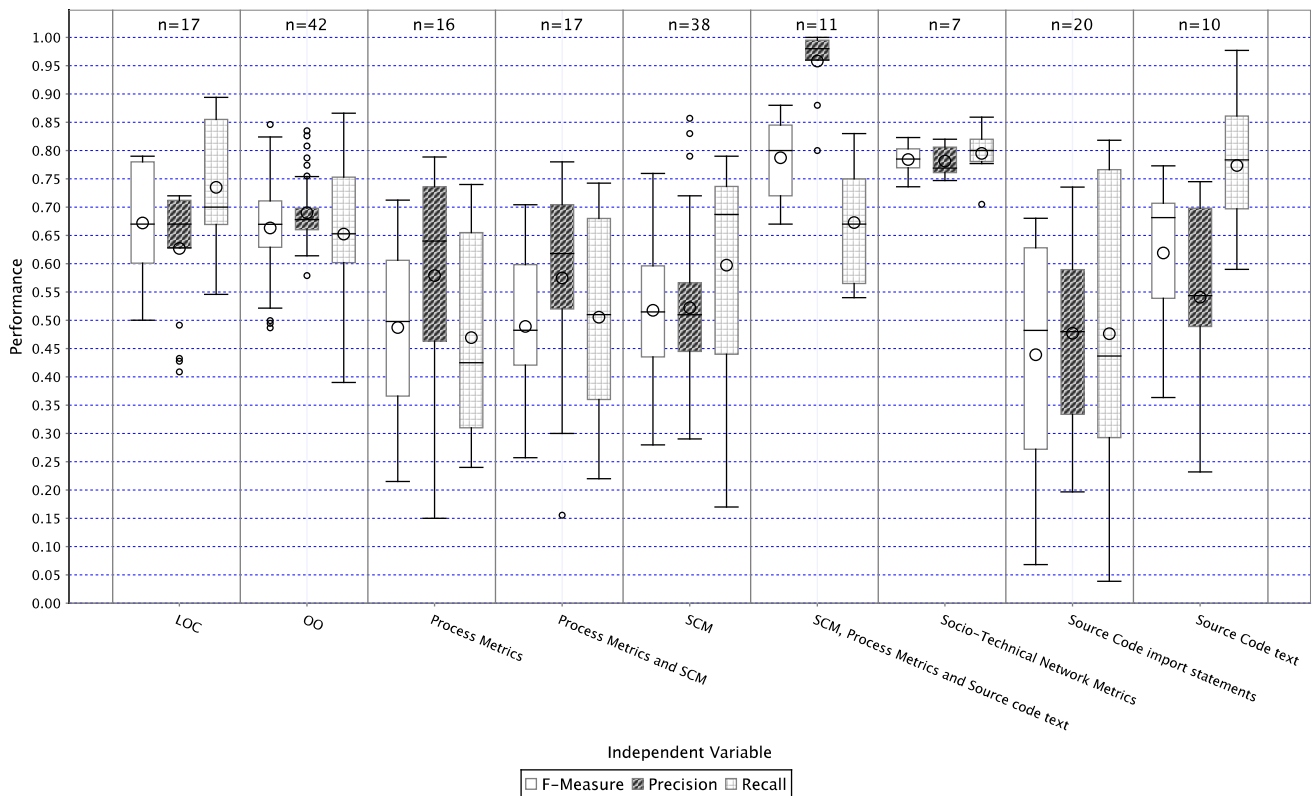


Fig. 7. Independent variables used in models.

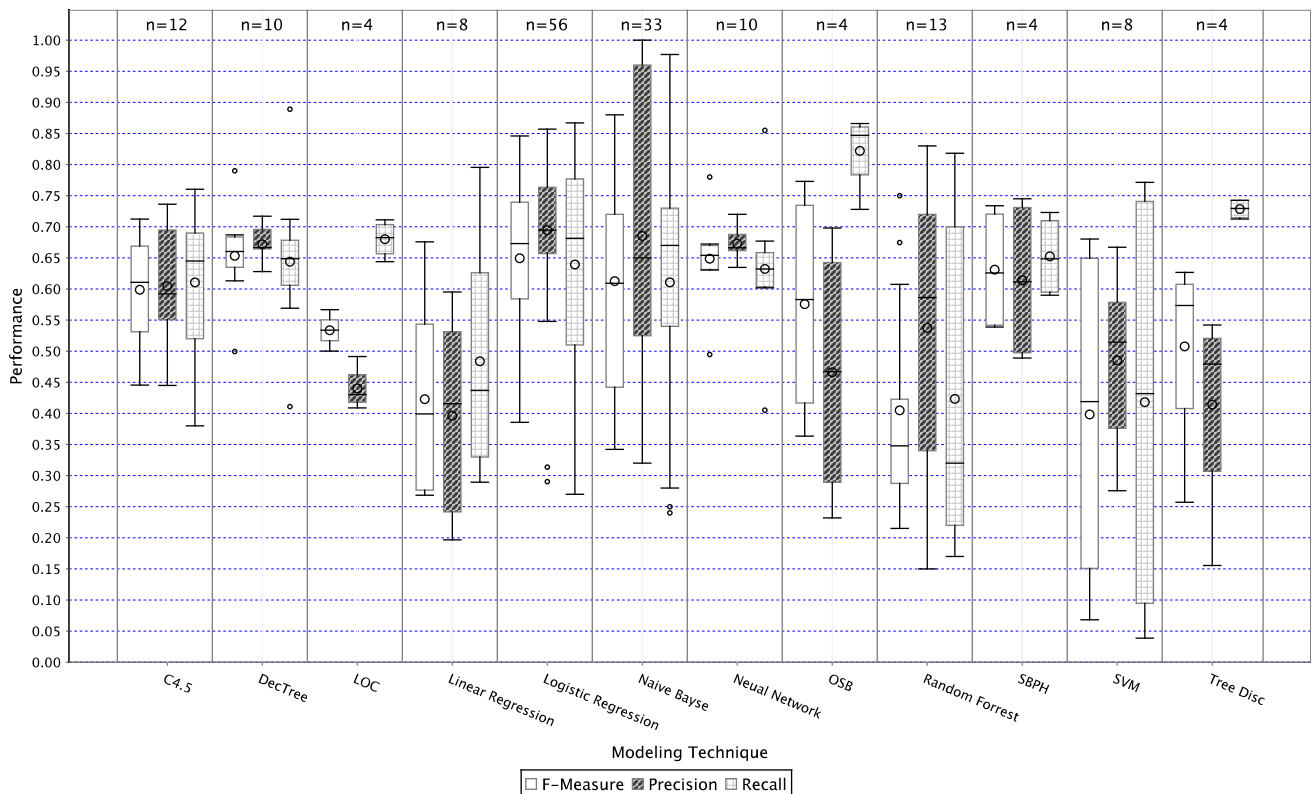


Fig. 8. Modeling technique used.

6 SYNTHESIS OF RESULTS

This section answers our research questions by synthesizing the qualitative and quantitative data we have collected. The qualitative data consist of the main findings reported by

each of the individual 36 finally included studies (presented in the Qualitative Data Table in our online Appendix). The quantitative data consist of the predictive performance of the individual models reported in the 36 studies (summarized in

the Categorical and Continuous Models Tables in our online Appendix). The quantitative data also consist of the detailed predictive performance data from 19 studies (206 models or model variants) comparing performance across models (reported in Section 5). This combination of data addresses model performance across studies and within individual studies. This allows us to discuss model performance in two ways. First, we discuss performance within individual studies to identify the main influences on model performance reported within a study. Second, we compare model performances across the models reported in 19 studies. This is an important approach to discussing fault prediction models. Most studies report at least one model which performs “well.” Though individual studies usually only compare performance within the set of models they present to identify their best model, we are able to then compare the performance of the models which perform well within a study across other studies. This allows us to report how well these models perform across studies.

6.1 Answering our Research Questions

RQ1: How does context affect fault prediction?

Analyzing model performance across the 19 studies in detail suggests that some context variables may influence the reliability of model prediction. Our results provide some evidence to suggest that predictive performance improves as systems get larger. This is suggested by the many models built for the Eclipse system. As Eclipse increases in size, the performance of models seems to improve. This makes some sense as models are likely to perform better with more data. We could find no evidence that this improved performance was based on the maturing of systems. It may be that size influences predictive performance more than system maturity. However, our dataset is relatively small and although we analyzed 206 models (or model variants) very few were based on immature systems. Our results also suggest that some applications may be less likely to produce reliable prediction models. For example, the many models built for embedded telecoms applications generally performed less well relative to other applications. Our results also show that many models have been built using Eclipse data. This corpus of knowledge on Eclipse provides a good opportunity for future researchers to meta-analyze across a controlled context.

The conventional wisdom is that context determines how transferrable a model is to other systems. Despite this, none of the 36 finally included studies directly investigate the impact on model performance of specific context variables such as system size, maturity, application area, or programming language. One exception is [S29], which demonstrates that transforming project data can make a model more comparable to other projects.

Many of the 36 finally included studies individually test how well their model performs when transferred to other contexts (releases, systems, application areas, data sources, or companies). Few of these studies directly investigate the contextual factors influencing the transferability of the model. Findings reported from individual studies on model transferability are varied. Most studies report that models perform poorly when transferred. In fact, Bell et al. [S11] report that models could not be

applied to other systems. Denaro and Pezzè [S37] reported good predictive performance only across homogenous applications. Nagappan et al. [S122] report that different subsets of complexity metrics relate to faults in different projects and that no single set of metrics fits all projects. Nagappan et al. [S122] conclude that models are only accurate when trained on the same or similar systems. However, other studies report more promising transferability. Weyuker et al. [S190] report good performance when models are transferred between releases of systems and between other systems. However, Shatnawi and Li [S160] report that the performance of models declines when applied to later releases of a system. Shatnawi and Li [S160] conclude that different metrics should be used in models used for later releases.

The context of models has not been studied extensively in the set of studies we analyzed. Although every model has been developed and tested within particular contexts, the impact of that context on model performance is scarcely studied directly. This is a significant gap in current knowledge as it means we currently do not know what context factors influence how well a model will transfer to other systems. It is therefore imperative that studies at least report their context since, in the future, this will enable a meta-analysis of the role context plays in predictive performance.

RQ2: Which independent variables should be included in fault prediction models?

Many different independent variables have been used in the 36 finally included studies. These mainly fall into process (e.g., previous change and fault data) and product (e.g., static code data) metrics as well as metrics relating to developers. In addition, some studies have used the text of the source code itself as the independent variables (e.g., Mizuno et al. [S116], Mizuno and Kikuno [S117]).

Model performance across the 19 studies we analyzed in detail suggests that the spam filtering technique, based on source code, used by Mizuno et al. [S116], Mizuno and Kikuno [S117] performs relatively well. On the other hand, models using only static code metrics (typically complexity-based) perform relatively poorly. Model performance does not seem to be improved by combining these metrics with OO metrics. Models seem to perform better using only OO metrics rather than only source code metrics. However, models using only LOC seem to perform just as well as those using only OO metrics and better than those models only using source code metrics. Within individual studies, Zhou et al. [S203] report that LOC data performs well. Ostrand et al. [S133] report that there was some value in LOC data and Hongyu [S56] reports LOC to be a useful early general indicator of fault-proneness. Zhou et al. [S203] report that LOC performs better than all but one of the Chidamber and Kemerer metrics (Weighted Methods per Class). Within other individual studies LOC data were reported to have poor predictive power and to be outperformed by other metrics (e.g., Bell et al. [S11]). Overall, LOC seem to be generally useful in fault prediction.

Model performance across the 19 studies that we analyzed suggests that the use of process data is not particularly related to good predictive performance. However, looking at

the findings from individual studies, several authors report that process data, in the form of previous history data, performs well (e.g., [S163], [S120]). D'Ambros et al. [S31] specifically report that previous bug reports are the best predictors. More sophisticated process measures have also been reported to perform well. In particular, Nagappan et al. [S120] introduce "change burst" metrics which demonstrate good predictive performance (however, these models perform only moderately when we compared them against models from other studies).

The few studies using developer information in models report conflicting results. Ostrand et al. [S135] report that the addition of developer information does not improve predictive performance much. Bird et al. [S18] report better performances when developer information is used as an element within a socio-technical network of variables. This study also performs well in our detailed comparison of performances (Bird et al. [S18] report results at a high level of granularity and so might be expected to perform better).

The models which perform best in our analysis of 19 studies seem to use a combined range of independent variables. For example, Shivaji et al. [S164] use process-based and SCM-based metrics together with source code. Bird et al. [S18] combine a range of metrics. The use of feature selection on sets of independent variables seems to improve the performance of models (e.g., [S164], [S76], [S18]). Optimized sets of metrics using, for example, feature selection, make sense.

RQ3: Which modeling techniques perform best when used in fault prediction?

While many included studies individually report the comparative performance of the modeling techniques they have used, no clear consensus on which perform best emerges when individual studies are looked at separately. Mizuno and Kikuno [S117] report that, of the techniques they studied, Orthogonal Sparse Bigrams Markov models (OSB) are best suited to fault prediction. Bibi et al. [S15] report that Regression via Classification (RvC) works well. Khoshgoftaar et al. [S86] report that modules whose fault proneness is predicted as uncertain can be effectively classified using the TreeDisc (TD) technique. Khoshgoftaar and Seliya [S83] also report that Case-Based Reasoning (CBR) does not predict well, with C4.5 also performing poorly. Arisholm et al. [S9] report that their comprehensive performance comparison revealed no predictive differences between the eight modeling techniques they investigated.

A clearer picture seems to emerge from our detailed analysis of model performance across the 19 studies. Our findings suggest that performance may actually be linked to the modeling technique used. Overall our comparative analysis suggests that studies using Support Vector Machine (SVM) techniques perform less well. These may be underperforming as they require parameter optimization (something rarely carried out in fault prediction studies) for best performance [18]. Where SVMs have been used in other prediction domains and may be better understood, they have performed well [19]. Models based on C4.5 seem to underperform if they use imbalanced data (e.g., Arisholm et al. [S8], [S9]), as the technique seems to be sensitive to this. Our comparative analysis also suggests that the models

performing comparatively well are relatively simple techniques that are easy to use and well understood. Naive Bayes and Logistic regression, in particular, seem to be the techniques used in models that are performing relatively well. Models seem to have performed best where the right technique has been selected for the right set of data. And these techniques have been tuned to the model (e.g., Shivaji et al. [S164]), rather than relying on default tool parameters.

7 METHODOLOGICAL ISSUES IN FAULT PREDICTION

The methodology used to develop, train, test, and measure the performance of fault prediction models is complex. However, the efficacy of the methodology used underpins the confidence which we can have in a model. It is essential that models use and report a rigorous methodology. Without this, the maturity of fault prediction in software engineering will be low. We identify methodological problems in existing studies so that future researchers can improve on these.

Throughout this SLR, methodological issues in the published studies came to light. During our assessment of the 208 initially included studies and the extraction of data from the 36 finally included studies methodological weaknesses emerged. In this section, we discuss the most significant of these methodological weaknesses. These generally relate to the quality of data used to build models and the approach taken to measure the predictive performance of models.

7.1 Data Quality

The quality of the data used in fault prediction has significant potential to undermine the efficacy of a model. Data quality is complex and many aspects of data are important to ensure reliable predictions. Unfortunately, it is often difficult to assess the quality of data used in studies, especially as many studies report very little about the data they use. Without good quality data, clearly reported, it is difficult to have confidence in the predictive results of studies.

The results of our assessment show that data quality is an issue in many studies. In fact many studies failed our synthesis assessment on the basis that they either reported insufficient information about the context of their data or about the collection of that data. Some studies explicitly acknowledge the importance of data quality (e.g., Jiang et al. [S64]).

Collecting good quality data is very hard. This is partly reflected by the number of studies which failed our assessment by not adequately explaining how they had collected their independent or dependent data. Fault data collection has been previously shown to be particularly hard to collect, usually because fault data are either not directly recorded or recorded poorly [20]. Collecting data is made more challenging because large datasets are usually necessary for reliable fault prediction. Jiang et al. [S64] investigate the impact that the size of the training and test dataset has on the accuracy of predictions. Tosun et al. [S176] present a useful insight into the real challenges associated with every aspect of fault prediction, but particularly on the difficulties of collecting reliable metrics and fault data. Once collected, data is usually noisy and

often needs to be cleaned (e.g., outliers and missing values dealt with [21]). Very few studies report any data cleaning (even in our 36 finally included studies).

The balance of data (i.e., the number of faulty as opposed to nonfaulty units) on which models are trained and tested is acknowledged by a few studies as fundamental to the reliability of models (see Appendix F for more information on class imbalance). Indeed, across the 19 studies we analyzed in detail, some of those performing best are based on data with a good proportion of faulty units (e.g., [S164], [S37], [S11], [S74]). Our analysis also suggests that data imbalance in relation to specific modeling techniques (e.g., C4.5) may be related to poor performance (e.g., [S8], [S9]). Several studies specifically investigated the impact of data balance and propose techniques to deal with it. For example, Khoshgoftaar et al. [S76] and Shivaji et al. [S164] present techniques for ensuring reliable data distributions. Schröter et al. [S154] base their training set on the faultiest parts of the system. Similarly, Seiffert et al. [S156] present data sampling and boosting techniques to address data imbalance. Data imbalance is explored further in Fioravanti and Nesi [S43] and Zhang et al. [S200]. Many studies seem to lack awareness of the need to account for data imbalance. As a consequence, the impact of imbalanced data on the real performance of models can be hidden by the performance measures selected. This is especially true where the balance of data is not even reported. Readers are then not able to account for the degree of imbalanced data in their interpretation of predictive performance.

7.2 Measuring the Predictive Performance of Models

There are many ways in which the performance of a prediction model can be measured. Indeed, many different categorical and continuous performance measures are used in our 36 studies. There is no one best way to measure the performance of a model. This depends on: the class distribution of the training data, how the model has been built, and how the model will be used. For example, the importance of measuring misclassification will vary depending on the application.

Performance comparison across studies is only possible if studies report a set of uniform measures. Furthermore, any uniform set of measures should give a full picture of correct and incorrect classification. To make models reporting categorical results most useful, we believe that the raw confusion matrix on which their performance measures are derived should be reported. This confusion matrix data would allow other researchers and potential users to calculate the majority of other measures. Pizzi et al. [22] provide a very usable format for presenting a confusion matrix. Some studies present many models and it is not practical to report the confusion matrices for all these. Menzies et al. [S114] suggest a useful way in which data from multiple confusion matrices may be effectively reported. Alternatively, Lessmann [S97] recommends that ROC curves and AUC are most useful when comparing the ability of modeling techniques to cope with different datasets (ROC curves do have some limitations [23]). Either of these approaches adopted widely would make studies more useful in the future. Comparing across studies reporting continuous

results is currently even more difficult and is the reason we were unable to present comparative boxplots across these studies. To enable cross comparison we recommend that continuous studies report Average Relative Error (ARE) in addition to any preferred measures presented.

The impact of performance measurement has been picked up in many studies. Zhou et al. [S203] report that the use of some measures, in the context of a particular model, can present a misleading picture of predictive performance and undermine the reliability of predictions. Arisholm et al. [S9] discuss how model performance varies depending on how it is measured. There is an increasing focus on identifying effective ways to measure the performance of models. Cost and/or effort aware measurement is now a significant strand of interest in prediction measurement. This takes into account the cost/effort of falsely identifying modules and has been increasingly reported as useful. The concept of cost-effectiveness measurement originated with the Simula group (e.g., Arisholm et al. [S9]), but has more recently been taken up and developed by other researchers, for example, Nagappan et al. [S120] and Mende and Koschke [S109].

7.3 Fault Severity

Few studies incorporate fault severity into their measurement of predictive performance. Although some faults are more important to identify than others, few models differentiate between the faults predicted. In fact, Shatnawi and Li's [S160] was the only study in the final 36 to use fault severity in their model. They report a model which is able to predict high and medium severity faults (these levels of severity are based on those reported in Bugzilla by Eclipse developers). Lamkanfi et al. [24], Singh et al. [S167], and Zhou and Leung [S202] are other studies which have also investigated severity. This lack of studies that consider severity is probably because, although acknowledged to be important, severity is considered a difficult concept to measure. For example, Menzies et al. [S113] say that severity is too vague to reliably investigate, Nikora and Munson [S126] says that "without a widely agreed definition of severity we cannot reason about it" and Ostrand et al. [S133] state that severity levels are highly subjective and can be inaccurate and inconsistent. These problems of how to measure and collect reliable severity data may limit the usefulness of fault prediction models. Companies developing noncritical systems may want to prioritize their fault finding effort only on the most severe faults.

7.4 The Reporting of Fault Prediction Studies

Our results suggest that, overall, fault prediction studies are reported poorly. Out of the 208 studies initially included in our review, only 36 passed our assessment criteria. Many of these criteria are focused on checking that studies report basic details about the study. Without a basic level of information reported it is hard to have confidence in a study. Our results suggest that many studies are failing to report information which is considered essential when reporting empirical studies in other domains. The poor reporting of studies has consequences for both future researchers and potential users of models: It is difficult for researchers to meta-analyze across studies and it is difficult

to replicate studies; it is also difficult for users to identify suitable models for implementation.

7.5 NASA Data

NASA's publicly available software metrics data have proven very popular in developing fault prediction models. We identify all 62 studies which use NASA data in the reference list of the 208 included studies. The NASA data is valuable as it enables studies using different modeling techniques and independent variables to be compared to others using the same dataset. It also allows studies to be replicated. A meta-analysis of the studies using NASA data would be valuable. However, although the repository holds many metrics and is publicly available, it does have limitations. It is not possible to explore the source code and the contextual data are not comprehensive (e.g., no data on maturity are available). It is also not always possible to identify if any changes have been made to the extraction and computation mechanisms over time. In addition, the data may suffer from important anomalies [21]. It is also questionable whether a model that works well on the NASA data will work on a different type of system; as Menzies et al. [S112] point out, NASA works in a unique niche market, developing software which is not typical of the generality of software systems. However, Turhan et al. [S181] have demonstrated that models built on NASA data are useful for predicting faults in software embedded in white goods.

8 THREATS TO VALIDITY

Searches. We do not include the term "quality" in our search terms as this would have resulted in the examination of a far wider range of irrelevant papers. This term generates a high number of false positive results. We might have missed some papers that use the term "quality" as a synonym for "defect" or "fault," etc. However, we missed only two papers that Catal and Diri's [2] searches found using the term "quality." This gives us confidence that we have missed very few papers. We also omitted the term "failure" from our search string as this generated papers predominately reporting on studies of software reliability in terms of safety critical systems. Such studies of reliability usually examine the dynamic behavior of the system and seldom look at the prediction of static code faults, which is the focus of this review.

We apply our search terms to only the titles of papers. We may miss studies that do not use these terms in the title. Since we extend our searches to include papers cited in the included papers, as well as key conferences, individual journals, and key authors, we are confident that the vast majority of key papers have been included.

Studies included for synthesis. The 36 studies which passed our assessment criteria may still have limitations that make their results unreliable. In the first place, the data on which these models are built might be problematic as we did not insist that studies report data cleaning or attribute selection. Nor did we apply any performance measure-based criteria. So some studies may be reporting unsafe predictive performances. This is a particular risk in regard to how studies have accounted for using imbalanced data.

This risk is mitigated in the categorical studies, where we are able to report precision, recall, and f-measure.

It is also possible that we have missed studies which should have been included in the set of 36 from which we extracted data. Some studies may have satisfied our assessment criteria but either failed to report what they did or did not report it in sufficient detail for us to be confident that they should pass the criteria. Similarly, we may have missed the reporting of a detail and a paper that should have passed a criterion did not. These risks are mitigated by two authors independently assessing every study.

The boxplots. The boxplots we present set performance against individual model factors (e.g., modeling technique used). This is a simplistic analysis, as a number of interacting factors are likely to underpin the performance of a model. For example, the technique used in combination with the dataset and the independent variables is likely to be more important than any one factor alone. Furthermore, methodological issues are also likely to impact on performance; for example, whether feature selection has been used. Our boxplots only present possible indicators of factors that should be investigated within the overall context of a model. More sophisticated analysis of a larger dataset is needed to investigate factors influencing model performance.

Our boxplots do not indicate the direction of any relationship between model performance and particular model factors. For example, we do not investigate whether a particular modeling technique performs well because it was used in a good model or whether a model performs well because it used a particular modeling technique. This is also important work for the future. In addition, some studies contribute data from many models to one boxplot, whereas other studies contribute data from only one model. This may skew the results. We do not calculate the statistical significance of any differences observed in the boxplots. This is because the data contained within them are not normally distributed and the individual points represent averages from different sizes of population.

9 CONCLUSIONS

Fault prediction is an important topic in software engineering. Fault prediction models have the potential to improve the quality of systems and reduce the costs associated with delivering those systems. As a result of this, many fault prediction studies in software engineering have been published. Our analysis of 208 of these studies shows that the vast majority are less useful than they could be. Most studies report insufficient contextual and methodological information to enable full understanding of a model. This makes it difficult for potential model users to select a model to match their context and few models have transferred into industrial practice. It also makes it difficult for other researchers to meta-analyze across models to identify the influences on predictive performance. A great deal of effort has gone into models that are of limited use to either practitioners or researchers.

The set of criteria we present identify a set of essential contextual and methodological details that fault prediction studies should report. These go some way toward addressing the need identified by Myrtveit et al. [25] for "more

reliable research procedures before we can have confidence in the conclusions of comparative studies." Our criteria should be used by future fault prediction researchers. They should also be used by journal and conference reviewers. This would ensure that future studies are built reliably and reported comparably with other such reliable studies. Of the 208 studies we reviewed, only 36 satisfied our criteria and reported essential contextual and methodological details.

We analyzed these 36 studies to determine what impacts on model performance in terms of the context of models, the independent variables used by models, and the modeling techniques on which they were built. Our results suggest that models which perform well tend to be built in a context where the systems are large. We found no evidence that the maturity of systems or the language used is related to predictive performance. But we did find some evidence to suggest that some application domains (e.g., embedded systems) may be more difficult to build reliable prediction models for. The independent variables used by models performing well seem to be sets of metrics (e.g., combinations of process, product, and people-based metrics). We found evidence that where models use KLOC as their independent variable, they perform no worse than where only single sets of other static code metrics are used. In addition, models which perform well tend to use simple, easy to use modeling techniques like Naive Bayes or Logistic Regression. More complex modeling techniques, such as support vector machines, tend to be used by models which perform relatively less well.

The methodology used to build models seems to be influential to predictive performance. The models which performed well seemed to optimize three aspects of the model. First, the choice of data was optimized. In particular, successful models tend to be trained on large datasets which contain a relatively high proportion of faulty units. Second, the choice of independent variables was optimized. A large range of metrics were used on which feature selection was applied. Third, the modeling technique was optimized. The default parameters were adjusted to ensure that the technique would perform effectively on the data provided.

Overall we conclude that many good fault prediction studies have been reported in software engineering (e.g., the 36 which passed our assessment criteria). Some of these studies are of exceptional quality, for example, Shivaji et al. [S164]. However, there remain many open questions about how to build effective fault prediction models for software systems. We need more studies which are based on a reliable methodology and which consistently report the context in which models are built and the methodology used to build them. A larger set of such studies will enable reliable cross-study metaanalysis of model performance. It will also give practitioners the confidence to appropriately select and apply models to their systems. Without this increase in reliable models that are appropriately reported, fault prediction will continue to have limited impact on the quality and cost of industrial software systems.

REFERENCES FOR THE 208 INCLUDED SLR PAPERS [S1-S208]

References from this list are cited using the format [Sref#]). Each reference is followed by a code indicating the status of

the paper in terms of whether it passed (P) or failed (F) our assessment. An indication is also given as to the assessment phase a paper failed (1, 2, 3, 4, or 5). The use of NASA data by studies is also indicated (N). A paper (n) failing an assessment criterion in phase 2 which used NASA data would be coded: (Paper=n; Status=F, Phase=2, Data=N)

Please report possible problems with our assessment of these papers to: tracy.hall@brunel.ac.uk.

- [1] R. Abreu and R. Premraj, "How Developer Communication Frequency Relates to Bug Introducing Changes," *Proc. Joint Int'l and Ann. ERCIM Workshops Principles of Software Evolution and Software Evolution Workshops*, pp. 153-158, 2009. (Paper=1, Status=F, Phase=1).
- [2] W. Afzal, "Using Faults-Slip-Through Metric as a Predictor of Fault-Proneness," *Proc. 17th Asia Pacific Software Eng. Conf.*, pp. 414-422, 2010. (Paper=2, Status=F, Phase=2).
- [3] W. Afzal and R. Torkar, "A Comparative Evaluation of Using Genetic Programming for Predicting Fault Count Data," *Proc. Third Int'l Conf. Software Eng. Advances*, pp. 407-414, 2008. (Paper=3, Status=F, Phase=2).
- [4] W. Afzal, R. Torkar, and R. Feldt, "Prediction of Fault Count Data Using Genetic Programming," *Proc. IEEE Int'l Multitopic Conf.*, pp. 349-356, 2008. (Paper=4, Status=F, Phase=2).
- [5] S. Amasaki, Y. Takagi, O. Mizuno, and T. Kikuno, "A Bayesian Belief Network for Assessing the Likelihood of Fault Content," *Proc. 14th Int'l Symp. Software Reliability Eng.*, pp. 215-226, Nov. 2003. (Paper=5, Status=F, Phase=2).
- [6] C. Andersson and P. Runeson, "A Replicated Quantitative Analysis of Fault Distributions in Complex Software Systems," *IEEE Trans. Software Eng.*, vol. 33, no. 5, pp. 273-286, May 2007. (Paper=6, Status=F, Phase=1).
- [7] E. Arisholm and L. Briand, "Predicting Fault-Prone Components in a Java Legacy System," *Proc. ACM/IEEE Int'l Symp. Empirical Software Eng.*, pp. 8-17, 2006. (Paper=7, Status=F, Phase=4).
- [8] E. Arisholm, L.C. Briand, and M. Fuglerud, "Data Mining Techniques for Building Fault-Proneness Models in Telecom Java Software," *Proc. IEEE 18th Int'l Symp. Software Reliability*, pp. 215-224, Nov. 2007. (Paper=8, Status=P).
- [9] E. Arisholm, L.C. Briand, and E.B. Johannessen, "A Systematic and Comprehensive Investigation of Methods to Build and Evaluate Fault Prediction Models," *J. Systems and Software*, vol. 83, no. 1, pp. 2-17, 2010. (Paper=9, Status=P).
- [10] N. Ayewah, W. Pugh, J. Morgenthaler, J. Penix, and Y. Zhou, "Evaluating Static Analysis Defect Warnings on Production Software," *Proc. Seventh ACM SIGPLAN-SIGSOFT Workshop Program Analysis for Software Tools and Eng.*, pp. 1-8, 2007. (Paper=10, Status=F, Phase=1).
- [11] R. Bell, T. Ostrand, and E. Weyuker, "Looking for Bugs in All the Right Places," *Proc. Int'l Symp. Software Testing and Analysis*, pp. 61-72, 2006. (Paper=11, Status=P).
- [12] P. Bellini, I. Bruno, P. Nesi, and D. Rogai, "Comparing Fault-proneness Estimation Models," *Proc. IEEE 10th Int'l Conf. Eng. Complex Computer Systems*, pp. 205-214, June 2005. (Paper=12, Status=F, Phase=2).
- [13] A. Bernstein, J. Ekanayake, and M. Pinzger, "Improving Defect Prediction Using Temporal Features and Non Linear Models," *Proc. Ninth Int'l Workshop Principles of Software Evolution: In Conjunction with the Sixth ESEC/FSE Joint Meeting*, pp. 11-18, 2007. (Paper=13, Status=F, Phase=2).
- [14] M. Bezerra, A. Oliveira, and S. Meira, "A Constructive RBF Neural Network for Estimating the Probability of Defects in Software Modules," *Proc. Int'l Joint Conf. Neural Networks*, pp. 2869-2874, Aug. 2007. (Paper=14, Status=F, Phase=2, Data=N).
- [15] S. Bibi, G. Tsoumakas, I. Stamelos, and I. Vlahvas, "Software Defect Prediction Using Regression via Classification," *Proc. IEEE Int'l Conf. Computer Systems and Applications*, vol. 8, pp. 330-336, 2006. (Paper=15, Status=P).
- [16] D. Binkley, H. Feild, D. Lawrie, and M. Pighin, "Software Fault Prediction Using Language Processing," *Proc. Testing: Academic and Industrial Conf. Practice and Research Techniques*, pp. 99-110, Sept. 2007. (Paper=16, Status=F, Phase=1).
- [17] D. Binkley, H. Feild, D. Lawrie, and M. Pighin, "Increasing Diversity: Natural Language Measures for Software Fault Prediction," *J. Systems and Software*, vol. 82, no. 11, pp. 1793-1803, 2009. (Paper=17, Status=F, Phase=1).

- [18] C. Bird, N. Nagappan, H. Gall, B. Murphy, and P. Devanbu, "Putting it All Together: Using Socio-Technical Networks to Predict Failures," *Proc. 20th Int'l Symp. Software Reliability Eng.*, pp. 109-119, 2009. (Paper=18, Status=P).
- [19] G. Boetticher, "Improving Credibility of Machine Learner Models in Software Engineering," *Advanced Machine Learner Applications in Software Eng.*, pp. 52-72, Idea Group Publishing, 2006. (Paper=19, Status=F, Phase=2, Data=N).
- [20] L. Briand, W. Melo, and J. Wust, "Assessing the Applicability of Fault-Proneness Models across Object-Oriented Software Projects," *IEEE Trans. Software Eng.*, vol. 28, no. 7, pp. 706-720, July 2002. (Paper=20, Status=F, Phase=2).
- [21] B. Caglayan, A. Bener, and S. Koch, "Merits of Using Repository Metrics in Defect Prediction for Open Source Projects," *Proc. ICSE Workshop Emerging Trends in Free/Libre/Open Source Software Research and Development*, pp. 31-36, 2009. (Paper=21, Status=P).
- [22] G. Calikli, A. Tosun, A. Bener, and M. Celik, "The Effect of Granularity Level on Software Defect Prediction," *Proc. 24th Int'l Symp. Computer and Information Sciences*, pp. 531-536, Sept. 2009. (Paper=22, Status=F, Phase=2).
- [23] C. Catal, B. Diri, and B. Ozumut, "An Artificial Immune System Approach for Fault Prediction in Object-Oriented Software," *Proc. Second Int'l Conf. Dependability of Computer Systems*, pp. 238-245, June 2007. (Paper=23, Status=F, Phase=2, Data=N).
- [24] E. Ceylan, F. Kutlubay, and A. Bener, "Software Defect Identification Using Machine Learning Techniques," *Proc. 32nd Software Eng. and Advanced Applications*, pp. 240-247, 2006. (Paper=24, Status=F, Phase=2).
- [25] V.U. Challagulla, F.B. Bastani, and I.-L. Yen, "A Unified Framework for Defect Data Analysis Using the MBR Technique," *Proc. IEEE 18th Int'l Tools with Artificial Intelligence*, pp. 39-46, Nov. 2006. (Paper=25, Status=F, Phase=2, Data=N).
- [26] V. Challagulla, F. Bastani, I.-L. Yen, and R. Paul, "Empirical Assessment of Machine Learning Based Software Defect Prediction Techniques," *Proc. S 10th IEEE Int'l Workshop Object-Oriented Real-Time Dependable Systems*, pp. 263-270, Feb. 2005. (Paper=26, Status=F, Phase=2, Data=N).
- [27] J. Cong, D. En-Mei, and Q. Li-Na, "Software Fault Prediction Model Based on Adaptive Dynamical and Median Particle Swarm Optimization," *Proc. Second Int'l Conf. Multimedia and Information Technology*, vol. 1, pp. 44-47, 2010. (Paper=27, Status=F, Phase=2, Data=N).
- [28] C. Cruz and A. Erika, "Exploratory Study of a UML Metric for Fault Prediction," *Proc. ACM/IEEE 32nd Int'l Conf. Software Eng.*, pp. 361-364, 2010. (Paper=28, Status=F, Phase=4).
- [29] C. Cruz, A. Erika, and O. Koichiro, "Towards Logistic Regression Models for Predicting Fault-Prone Code across Software Projects," *Proc. Third Int'l Symp. Empirical Software Eng. and Measurement*, pp. 460-463, 2009. (Paper=29, Status=P).
- [30] V. Dallmeier and T. Zimmermann, "Extraction of Bug Localization Benchmarks from History," *Proc. IEEE/ACM 22nd Int'l Conf. Automated Software Eng.*, pp. 433-436, 2007. (Paper=30, Status=F, Phase=1).
- [31] M. D'Ambros, M. Lanza, and R. Robbes, "On the Relationship between Change Coupling and Software Defects," *Proc. 16th Working Conf. Reverse Eng.*, pp. 135-144, Oct. 2009. (Paper=31, Status=P).
- [32] M. D'Ambros, M. Lanza, and R. Robbes, "An Extensive Comparison of Bug Prediction Approaches," *Proc. IEEE Seventh Working Conf. Mining Software Repositories*, pp. 31-41, 2010. (Paper=32, Status=P).
- [33] A.B. de Carvalho, A. Pozo, S. Vergilio, and A. Lenz, "Predicting Fault Proneness of Classes through a Multiobjective Particle Swarm Optimization Algorithm," *Proc. IEEE 20th Int'l Conf. Tools with Artificial Intelligence*, vol. 2, pp. 387-394, 2008. (Paper=33, Status=F, Phase=2, Data=N).
- [34] A.B. de Carvalho, A. Pozo, and S.R. Vergilio, "A Symbolic Fault-Prediction Model Based on Multiobjective Particle Swarm Optimization," *J. Systems and Software*, vol. 83, no. 5, pp. 868-882, 2010. (Paper=34, Status=F, Phase=2, Data=N).
- [35] G. Denaro, "Estimating Software Fault-Proneness for Tuning Testing Activities," *Proc. Int'l Conf. Software Eng.*, pp. 704-706, 2000. (Paper=35, Status=F, Phase=2).
- [36] G. Denaro, S. Morasca, and M. Pezzè, "Deriving Models of Software Fault-Proneness," *Proc. 14th Int'l Conf. Software Eng. and Knowledge Eng.*, pp. 361-368, 2002. (Paper=36, Status=F, Phase=2).
- [37] G. Denaro and M. Pezzè, "An Empirical Evaluation of Fault-Proneness Models," *Proc. 24th Int'l Conf. Software Eng.*, pp. 241-251, 2002. (Paper=37, Status=P).
- [38] Z. Dianqin and W. Zhongyuan, "The Application of Gray-Prediction Theory in the Software Defects Management," *Proc. Int'l Conf. Computational Intelligence and Software Eng.*, pp. 1-5, 2009. (Paper=38, Status=F, Phase=2).
- [39] K. El Emam, W. Melo, and J. Machado, "The Prediction of Faulty Classes Using Object-Oriented Design Metrics," *J. Systems and Software*, vol. 56, no. 1, pp. 63-75, 2001. (Paper=39, Status=F, Phase=2).
- [40] K. Elish and M. Elish, "Predicting Defect-prone Software Modules Using Support Vector Machines," *J. Systems and Software*, vol. 81, no. 5, pp. 649-660, 2008. (Paper=40, Status=F, Phase=2, Data=N).
- [41] N. Fenton, M. Neil, W. Marsh, P. Hearty, L. Radlinski, and P. Krause, "Project Data Incorporating Qualitative Factors for Improved Software Defect Prediction," *Proc. Int'l Workshop Predictor Models in Software Eng.*, p. 2, May 2007. (Paper=41, Status=F, Phase=2).
- [42] N. Fenton and N. Ohlsson, "Quantitative Analysis of Faults and Failures in a Complex Software System," *IEEE Trans. Software Eng.*, vol. 26, no. 8, pp. 797-814, Aug. 2000. (Paper=42, Status=F, Phase=1).
- [43] F. Fioravanti and P. Nesi, "A Study on Fault-Proneness Detection of Object-Oriented Systems," *Proc. Fifth European Conf. Software Maintenance and Reeng.*, pp. 121-130, 2001. (Paper=43, Status=F, Phase=2).
- [44] K. Gao and T. Khoshgoftaar, "A Comprehensive Empirical Study of Count Models for Software Fault Prediction," *IEEE Trans. Reliability*, vol. 56, no. 2, pp. 223-236, June 2007. (Paper=44, Status=F, Phase=2).
- [45] N. Gayatri, S. Nickolas, A.V. Reddy, and R. Chitra, "Performance Analysis of Datamining Algorithms for Software Quality Prediction," *Proc. Int'l Conf. Advances in Recent Technologies in Comm. and Computing*, pp. 393-395, 2009. (Paper=45, Status=F, Phase=2, Data=N).
- [46] I. Gondra, "Applying Machine Learning to Software Fault-Proneness Prediction," *J. Systems and Software*, vol. 81, no. 2, pp. 186-195, 2008. (Paper=46, Status=F, Phase=1).
- [47] T. Graves, A. Karr, J. Marron, and H. Siy, "Predicting Fault Incidence Using Software Change History," *IEEE Trans. Software Eng.*, vol. 26, no. 7, pp. 653-661, July 2000. (Paper=47, Status=F, Phase=1).
- [48] D. Gray, D. Bowes, N. Davey, S. Yi, and B. Christianson, "Software Defect Prediction Using Static Code Metrics Underestimates Defect-Proneness," *Proc. Int'l Joint Conf. Neural Networks*, pp. 1-7, 2010. (Paper=48, Status=F, Phase=1).
- [49] L. Guo, B. Cukic, and H. Singh, "Predicting Fault Prone Modules by the Dempster-Shafer Belief Networks," *Proc. IEEE 18th Int'l Conf. Automated Software Eng.*, pp. 249-252, Oct. 2003. (Paper=49, Status=F, Phase=2, Data=N).
- [50] L. Guo, Y. Ma, B. Cukic, and H. Singh, "Robust Prediction of Fault-Proneness by Random Forests," *Proc. 15th Int'l Symp. Software Reliability Eng.*, pp. 417-428, Nov. 2004. (Paper=50, Status=F, Phase=2, Data=N).
- [51] T. Gyimothy, R. Ferenc, and I. Siket, "Empirical Validation of Object-Oriented Metrics on Open Source Software for Fault Prediction," *IEEE Trans. Software Eng.*, vol. 31, no. 10, pp. 897-910, Oct. 2005. (Paper=51, Status=P).
- [52] A.E. Hassan, "Predicting Faults Using the Complexity of Code Changes," *Proc. 31st IEEE Int'l Conf. Software Eng.*, pp. 78-88, 2009. (Paper=52, Status=F, Phase=5).
- [53] A. Hassan and R. Holt, "The Top Ten List: Dynamic Fault Prediction," *Proc. 21st IEEE Int'l Conf. Software Maintenance*, pp. 263-272, Sept. 2005. (Paper=53, Status=F, Phase=1).
- [54] Y. Higo, K. Murao, S. Kusumoto, and K. Inoue, "Predicting Fault-Prone Modules Based on Metrics Transitions," *Proc. Workshop Defects in Large Software Systems*, pp. 6-10, 2008. (Paper=54, Status=F, Phase=1).
- [55] T. Holschuh, M. Pauser, K. Herzig, T. Zimmermann, R. Premraj, and A. Zeller, "Predicting Defects in SAP Java Code: An Experience Report," *Proc. 31st Int'l Conf. Software Eng.-Companion Volume*, pp. 172-181, 2009. (Paper=55, Status=F, Phase=2).
- [56] Z. Hongyu, "An Investigation of the Relationships Between Lines of Code and Defects," *Proc. IEEE Int'l Conf. Software Maintenance*, pp. 274-283, 2009. (Paper=56, Status=P, Data=N).

- [57] D. Hovemeyer and W. Pugh, "Finding Bugs Is Easy," *Proc. Companion to the 19th Ann. ACM SIGPLAN Conf. Object-Oriented Programming Systems, Languages, and Applications*, pp. 132-136, 2004. (Paper=57, Status=F, Phase=3).
- [58] L. Hribar and D. Duka, "Software Component Quality Prediction Using KNN and Fuzzy Logic," *Proc. 33rd Int'l Convention MIPRO*, pp. 402-408, 2010. (Paper=58, Status=F, Phase=1).
- [59] W. Huanjing, T.M. Khoshgoftaar, and A. Napolitano, "A Comparative Study of Ensemble Feature Selection Techniques for Software Defect Prediction," *Proc. Ninth Int'l Conf. Machine Learning and Applications*, pp. 135-140, 2010. (Paper=59, Status=F, Phase=4, Data=N).
- [60] L. Jiang, Z. Su, and E. Chiu, "Context-Based Detection of Clone-Related Bugs," *Proc. Sixth Joint Meeting of the European Software Eng. Conf. and the ACM SIGSOFT Symp. Foundations of Software Eng.*, pp. 55-64, 2007. (Paper=60, Status=F, Phase=1).
- [61] Y. Jiang, B. Cukic, and Y. Ma, "Techniques for Evaluating Fault Prediction Models," *Empirical Software Eng.*, vol. 13, no. 5, pp. 561-595, 2008. (Paper=61, Status=F, Phase=2, Data=N).
- [62] Y. Jiang, B. Cukic, and T. Menzies, "Fault Prediction Using Early Lifecycle Data," *Proc. IEEE 18th Int'l Symp. Software Reliability*, pp. 237-246, Nov. 2007. (Paper=62, Status=F, Phase=2, Data=N).
- [63] Y. Jiang, B. Cukic, and T. Menzies, "Cost Curve Evaluation of Fault Prediction Models," *Proc. 19th Int'l Symp. Software Reliability Eng.*, pp. 197-206, Nov. 2008. (Paper=63, Status=F, Phase=2, Data=N).
- [64] Y. Jiang, J. Lin, B. Cukic, and T. Menzies, "Variance Analysis in Software Fault Prediction Models," *Proc. 20th Int'l Symp. Software Reliability Eng.*, pp. 99-108, Nov. 2009. (Paper=64, Status=F, Phase=2, Data=N).
- [65] J.A. Jones and M.J. Harrold, "Empirical Evaluation of the Tarantula Automatic Fault-Localization Technique," *Proc. IEEE/ACM 20th Int'l Conf. Automated Software Eng.*, pp. 273-282, 2005. (Paper=65, Status=F, Phase=1).
- [66] J.A. Jones, M.J. Harrold, and J. Stasko, "Visualization of Test Information to Assist Fault Localization," *Proc. 24th Int'l Conf. Software Eng.*, pp. 467-477, 2002. (Paper=66, Status=F, Phase=1).
- [67] H. Joshi, C. Zhang, S. Ramaswamy, and C. Bayrak, "Local and Global Recency Weighting Approach to Bug Prediction," *Proc. Fourth Int'l Workshop Mining Software Repositories*, p. 33, May 2007. (Paper=67, Status=F, Phase=2).
- [68] L. Jun, X. Zheng, Q. Jianzhong, and L. Shukuan, "A Defect Prediction Model for Software Based on Service Oriented Architecture Using Expert Cocomo," *Proc. 21st Ann. Int'l Conf. Chinese Control and Decision Conf.*, pp. 2591-2594, 2009. (Paper=68, Status=F, Phase=1).
- [69] Y. Kamei, S. Matsumoto, A. Monden, K.i. Matsumoto, B. Adams, and A.E. Hassan, "Revisiting Common Bug Prediction Findings Using Effort-Aware Models," *Proc. IEEE Int'l Conf. Software Maintenance*, pp. 1-10, 2010. (Paper=69, Status=P).
- [70] K. Kaminsky and G. Boetticher, "Building a Genetically Engineerable Evolvable Program (GEEP) Using Breadth-based Explicit Knowledge for Predicting Software Defects," *Proc. IEEE Ann. Meeting Fuzzy Information*, vol. 1, pp. 10-15, June 2004. (Paper=70, Status=F, Phase=1).
- [71] S. Kanmani, V.R. Uthariaraj, V. Sankaranarayanan, and P. Thambidurai, "Object Oriented Software Quality Prediction Using General Regression Neural Networks," *SIGSOFT Software Eng. Notes*, vol. 29, pp. 1-6, Sept. 2004. (Paper=71, Status=F, Phase=2).
- [72] S. Kanmani, V. Uthariaraj, V. Sankaranarayanan, and P. Thambidurai, "Object-Oriented Software Fault Prediction Using Neural Networks," *Information and Software Technology*, vol. 49, no. 5, pp. 483-492, 2007. (Paper=72, Status=F, Phase=2).
- [73] Y. Kastro and A. Bener, "A Defect Prediction Method for Software Versioning," *Software Quality J.*, vol. 16, no. 4, pp. 543-562, 2008. (Paper=73, Status=P).
- [74] A. Kaur and R. Malhotra, "Application of Random Forest in Predicting Fault-Prone Classes," *Proc. Int'l Conf. Advanced Computer Theory and Eng.*, pp. 37-43, 2008. (Paper=74, Status=P).
- [75] A. Kaur, P.S. Sandhu, and A.S. Bra, "Early Software Fault Prediction Using Real Time Defect Data," *Proc. Second Int'l Conf. Machine Vision*, pp. 242-245, 2009. (Paper=75, Status=F, Phase=2, Data=N).
- [76] T.M. Khoshgoftaar, K. Gao, and N. Seliya, "Attribute Selection and Imbalanced Data: Problems in Software Defect Prediction," *Proc. 22nd IEEE Int'l Conf. Tools with Artificial Intelligence*, vol. 1, pp. 137-144, 2010. (Paper=76, Status=P).
- [77] T. Khoshgoftaar, E. Allen, and J. Busboom, "Modeling Software Quality: The Software Measurement Analysis and Reliability Toolkit," *Proc. IEEE 12th Int'l Conf. Tools with Artificial Intelligence*, pp. 54-61, 2000. (Paper=77, Status=F, Phase=2).
- [78] T. Khoshgoftaar, K. Gao, and R. Szabo, "An Application of Zero-Inflated Poisson Regression for Software Fault Prediction," *Proc. 12th Int'l Symp. Software Reliability Eng.*, pp. 66-73, Nov. 2001. (Paper=78, Status=F, Phase=2).
- [79] T. Khoshgoftaar, E. Geleyn, and K. Gao, "An Empirical Study of the Impact of Count Models Predictions on Module-order Models," *Proc. Eighth IEEE Symp. Software Metrics*, pp. 161-172, 2002. (Paper=79, Status=F, Phase=2).
- [80] T. Khoshgoftaar and N. Seliya, "Improving Usefulness of Software Quality Classification Models Based on Boolean Discriminant Functions," *Proc. 13th Int'l Symp. Software Reliability Eng.*, pp. 221-230, 2002. (Paper=80, Status=F, Phase=2).
- [81] T. Khoshgoftaar and N. Seliya, "Software Quality Classification Modeling Using the Sprint Decision Tree Algorithm," *Proc. IEEE 14th Int'l Conf. Tools with Artificial Intelligence*, pp. 365-374, 2002. (Paper=81, Status=F, Phase=2).
- [82] T. Khoshgoftaar and N. Seliya, "Tree-Based Software Quality Estimation Models for Fault Prediction," *Proc. IEEE Eighth Symp. Software Metrics*, pp. 203-214, 2002. (Paper=82, Status=F, Phase=2).
- [83] T. Khoshgoftaar and N. Seliya, "Comparative Assessment of Software Quality Classification Techniques: An Empirical Case Study," *Empirical Software Eng.*, vol. 9, no. 3, pp. 229-257, 2004. (Paper=83, Status=P).
- [84] T. Khoshgoftaar, N. Seliya, and K. Gao, "Assessment of a New Three-Group Software Quality Classification Technique: An Empirical Case Study," *Empirical Software Eng.*, vol. 10, no. 2, pp. 183-218, 2005. (Paper=84, Status=F, Phase=2).
- [85] T. Khoshgoftaar, V. Thaker, and E. Allen, "Modeling Fault-Prone Modules of Subsystems," *Proc. 11th Int'l Symp. Software Reliability Eng.*, pp. 259-267, 2000. (Paper=85, Status=F, Phase=2).
- [86] T. Khoshgoftaar, X. Yuan, E. Allen, W. Jones, and J. Hudepohl, "Uncertain Classification of Fault-Prone Software Modules," *Empirical Software Eng.*, vol. 7, no. 4, pp. 297-318, 2002. (Paper=86, Status=P).
- [87] S. Kim, K. Pan, and E. Whitehead Jr., "Memories of Bug Fixes," *Proc. 14th ACM SIGSOFT Int'l Symp. Foundations of Software Eng.*, pp. 35-45, 2006. (Paper=87, Status=F, Phase=2).
- [88] S. Kim, T. Zimmermann, E. Whitehead Jr., and A. Zeller, "Predicting Faults from Cached History," *Proc. 29th Int'l Conf. Software Eng.*, pp. 489-498, 2007. (Paper=88, Status=F, Phase=4).
- [89] S. Kim, T. Zimmermann, K. Pan, and E. Whitehead, "Automatic Identification of Bug-Introducing Changes," *Proc. IEEE/ACM 21st Int'l Conf. Automated Software Eng.*, pp. 81-90, Sept. 2006. (Paper=89, Status=F, Phase=1).
- [90] M. Kläs, F. Elberzhager, J. Münch, K. Hartjes, and O. von Graevemeyer, "Transparent Combination of Expert and Measurement Data for Defect Prediction: An Industrial Case Study," *Proc. 32nd ACM/IEEE Int'l Conf. Software Eng.*, pp. 119-128, 2010. (Paper=90, Status=F, Phase=2).
- [91] M. Kläs, H. Nakao, F. Elberzhager, and J. Münch, "Predicting Defect Content and Quality Assurance Effectiveness by Combining Expert Judgment and Defect Data-A Case Study," *Proc. 19th Int'l Symp. Software Reliability Eng.*, pp. 17-26, 2008. (Paper=91, Status=F, Phase=2).
- [92] P. Knab, M. Pinzger, and A. Bernstein, "Predicting Defect Densities in Source Code Files with Decision Tree Learners," *Proc. Int'l Workshop Mining Software Repositories*, pp. 119-125, 2006. (Paper=92, Status=P).
- [93] A. Koru and H. Liu, "Building Effective Defect-Prediction Models in Practice," *IEEE Software*, vol. 22, no. 6, pp. 23-29, Nov./Dec. 2005. (Paper=93, Status=F, Phase=2, Data=N).
- [94] A. Koru, D. Zhang, and H. Liu, "Modeling the Effect of Size on Defect Proneness for Open-Source Software," *Proc. Int'l Workshop Predictor Models in Software Eng.*, p. 10, May 2007. (Paper=94, Status=F, Phase=1).
- [95] O. Kutlubay, B. Turhan, and A. Bener, "A Two-Step Model for Defect Density Estimation," *Proc. 33rd EUROMICRO Conf. Software Eng. and Advanced Applications*, pp. 322-332, Aug. 2007. (Paper=95, Status=F, Phase=2, Data=N).

- [96] L. Layman, G. Kudrjavets, and N. Nagappan, "Iterative Identification of Fault-Prone Binaries Using in-Process Metrics," *Proc. ACM-IEEE Second Int'l Symp. Empirical Software Eng. and Measurement*, pp. 206-212, 2008. (Paper=96, Status=F, Phase=4).
- [97] S. Lessmann, B. Baesens, C. Mues, and S. Pietsch, "Benchmarking Classification Models for Software Defect Prediction: A Proposed Framework and Novel Findings," *IEEE Trans. Software Eng.*, vol. 34, no. 4, pp. 485-496, July/Aug. 2008. (Paper=97, Status=F, Phase=2, Data=N).
- [98] P.L. Li, J. Herbsleb, M. Shaw, and B. Robinson, "Experiences and Results from Initiating Field Defect Prediction and Product Test Prioritization Efforts at Abb Inc." *Proc. 28th Int'l Conf. Software Eng.*, pp. 413-422, 2006. (Paper=98, Status=P).
- [99] P.L. Li, M. Shaw, J. Herbsleb, B. Ray, and P. Santhanam, "Empirical Evaluation of Defect Projection Models for Widely-Deployed Production Software Systems," *SIGSOFT Software Eng. Notes*, vol. 29, pp. 263-272, Oct. 2004. (Paper=99, Status=F, Phase=1).
- [100] P. Li, J. Herbsleb, and M. Shaw, "Finding Predictors of Field Defects for Open Source Software Systems in Commonly Available Data Sources: A Case Study of Opensbd," *Proc. IEEE 11th Int'l Symp. Software Metrics*, pp. 10-32, Sept. 2005. (Paper=100, Status=F, Phase=1).
- [101] Z. Li and M. Reformat, "A Practical Method for the Software Fault-Prediction," *Proc. IEEE Int'l Conf. Information Reuse and Integration*, pp. 659-666, Aug. 2007. (Paper=101, Status=F, Phase=2, Data=N).
- [102] Y. Ma, L. Guo, and B. Cukic, "A Statistical Framework for the Prediction of Fault-Proneness," *Proc. Advances in Machine Learning Applications in Software Eng.*, pp. 237-265, 2006. (Paper=102, Status=F, Phase=2, Data=N).
- [103] J.T. Madhavan and E.J. Whitehead Jr., "Predicting Buggy Changes Inside an Integrated Development Environment," *Proc. OOPSLA Workshop Eclipse Technology Exchange*, pp. 36-40, 2007. (Paper=103, Status=F, Phase=4).
- [104] A. Mahaweerawat, P. Sophatsathit, and C. Lursinsap, "Software Fault Prediction Using Fuzzy Clustering and Radial-Basis Function Network," *Proc. Int'l Conf. Intelligent Technologies*, pp. 304-313, 2002. (Paper=104, Status=F, Phase=2).
- [105] A. Mahaweerawat, P. Sophatsathit, and C. Lursinsap, "Adaptive Self-Organizing Map Clustering for Software Fault Prediction," *Proc. Fourth Int'l Joint Conf. Computer Science and Software Eng.*, pp. 35-41, 2007. (Paper=105, Status=F, Phase=2).
- [106] A. Mahaweerawat, P. Sophatsathit, C. Lursinsap, and P. Musilek, "Fault Prediction in Object-Oriented Software Using Neural Network Techniques," *Proc. InTech Conf.*, pp. 2-4, 2004. (Paper=106, Status=F, Phase=2).
- [107] A. Marcus, D. Poshyanyk, and R. Ferenc, "Using the Conceptual Cohesion of Classes for Fault Prediction in Object-Oriented Systems," *IEEE Trans. Software Eng.*, vol. 34, no. 2, pp. 287-300, Mar./Apr. 2008. (Paper=107, Status=F, Phase=1).
- [108] T. Mende and R. Koschke, "Revisiting the Evaluation of Defect Prediction Models," *Proc. Fifth Int'l Conf. Predictor Models in Software Eng.*, p. 7, 2009. (Paper=108, Status=F, Phase=2, Data=N).
- [109] T. Mende and R. Koschke, "Effort-Aware Defect Prediction Models," *Proc. 14th European Conf. Software Maintenance and Reeng.*, pp. 107-116, 2010. (Paper=109, Status=P, Data=N).
- [110] T. Mende, R. Koschke, and M. Leszak, "Evaluating Defect Prediction Models for a Large Evolving Software System," *Proc. 13th European Conf. Software Maintenance and Reeng.*, pp. 247-250, Mar. 2009. (Paper=110, Status=P).
- [111] T. Menzies and J. Di Stefano, "How Good Is Your Blind Spot Sampling Policy," *Proc. IEEE Eighth Int'l Symp. High Assurance Systems Eng.*, pp. 129-138, Mar. 2004. (Paper=111, Status=F, Phase=2, Data=N).
- [112] T. Menzies, J. Greenwald, and A. Frank, "Data Mining Static Code Attributes to Learn Defect Predictors," *IEEE Trans. Software Eng.*, vol. 33, no. 1, pp. 2-13, Jan. 2007. (Paper=112, Status=F, Phase=2, Data=N).
- [113] T. Menzies, Z. Milton, B. Turhan, B. Cukic, Y. Jiang, and A.B. Bener, "Defect Prediction from Static Code Features: Current Results, Limitations, New Approaches," *Automatic Software Eng.*, vol. 17, no. 4, pp. 375-407, 2010. (Paper=113, Status=F, Phase=2, Data=N).
- [114] T. Menzies, B. Turhan, A. Bener, G. Gay, B. Cukic, and Y. Jiang, "Implications of Ceiling Effects in Defect Predictors," *Proc. Fourth Int'l Workshop Predictor Models in Software Eng.*, pp. 47-54, 2008. (Paper=114, Status=F, Phase=2, Data=N).
- [115] M. Mertik, M. Lenic, G. Stiglic, and P. Kokol, "Estimating Software Quality with Advanced Data Mining Techniques," *Proc. Int'l Conf. Software Eng. Advances*, p. 19, Oct. 2006. (Paper=115, Status=F, Phase=2, Data=N).
- [116] O. Mizuno, S. Ikami, S. Nakaichi, and T. Kikuno, "Spam Filter Based Approach for Finding Fault-Prone Software Modules," *Proc. Fourth Int'l Workshop Mining Software Repositories*, p. 4, May 2007. (Paper=116, Status=P).
- [117] O. Mizuno and T. Kikuno, "Training on Errors Experiment to Detect Fault-Prone Software Modules by Spam Filter," *Proc. Sixth Joint Meeting of the European Software Eng. Conf. and the ACM SIGSOFT Symp. Foundations of Software Eng.*, pp. 405-414, 2007. (Paper=117, Status=P).
- [118] R. Moser, W. Pedrycz, and G. Succi, "A Comparative Analysis of the Efficiency of Change Metrics and Static Code Attributes for Defect Prediction," *Proc. ACM/IEEE 30th Int'l Conf. Software Eng.*, pp. 181-190, 2008. (Paper=118, Status=P).
- [119] N. Nagappan and T. Ball, "Static Analysis Tools as Early Indicators of Pre-Release Defect Density," *Proc. 27th Int'l Conf. Software Eng.*, pp. 580-586, May 2005. (Paper=119, Status=F, Phase=3).
- [120] N. Nagappan, A. Zeller, T. Zimmermann, K. Herzig, and B. Murphy, "Change Bursts as Defect Predictors," *Proc. IEEE 21st Int'l Symp. Software Reliability Eng.*, pp. 309-318, 2010. (Paper=120, Status=P).
- [121] N. Nagappan and T. Ball, "Use of Relative Code Churn Measures to Predict System Defect Density," *Proc. 27th Int'l Conf. Software Eng.*, pp. 284-292, 2005. (Paper=121, Status=F, Phase=2).
- [122] N. Nagappan, T. Ball, and A. Zeller, "Mining Metrics to Predict Component Failures," *Proc. 28th Int'l Conf. Software Eng.*, pp. 452-461, 2006. (Paper=122, Status=P).
- [123] N.K. Nagwani and S. Verma, "Predictive Data Mining Model for Software Bug Estimation Using Average Weighted Similarity," *Proc. IEEE Second Int'l Advance Computing Conf.*, pp. 373-378, 2010. (Paper=123, Status=F, Phase=1).
- [124] A. Neufelder, "How to Measure the Impact of Specific Development Practices on Fielded Defect Density," *Proc. 11th Int'l Symp. Software Reliability Eng.*, pp. 148-160, 2000. (Paper=124, Status=F, Phase=1).
- [125] A. Nikora and J. Munson, "Developing Fault Predictors for Evolving Software Systems," *Proc. Ninth Int'l Software Metrics Symp.*, pp. 338-350, Sept. 2003. (Paper=125, Status=F, Phase=1).
- [126] A. Nikora and J. Munson, "The Effects of Fault Counting Methods on Fault Model Quality," *Proc. 28th Ann. Int'l Computer Software and Applications Conf.*, vol. 1, pp. 192-201, Sept. 2004. (Paper=126, Status=F, Phase=1).
- [127] A. Nugroho, M.R.V. Chaudron, and E. Arisholm, "Assessing UML Design Metrics for Predicting Fault-Prone Classes in a Java System," *Proc. IEEE Seventh Working Conf. Mining Software Repositories*, pp. 21-30, 2010. (Paper=127, Status=F, Phase=5).
- [128] H. Olague, L. Etzkorn, S. Gholston, and S. Quattlebaum, "Empirical Validation of Three Software Metrics Suites to Predict Fault-Proneness of Object-Oriented Classes Developed Using Highly Iterative or Agile Software Development Processes," *IEEE Trans. Software Eng.*, vol. 33, no. 6, pp. 402-419, June 2007. (Paper=128, Status=F, Phase=2).
- [129] A. Oral and A. Bener, "Defect Prediction for Embedded Software," *Proc. 22nd Int'l Symp. Computer and Information Sciences*, pp. 1-6, Nov. 2007. (Paper=129, Status=F, Phase=2, Data=N).
- [130] T.J. Ostrand and E.J. Weyuker, "The Distribution of Faults in a Large Industrial Software System," *SIGSOFT Software Eng. Notes*, vol. 27, pp. 55-64, July 2002. (Paper=130, Status=F, Phase=1).
- [131] T.J. Ostrand, E.J. Weyuker, and R.M. Bell, "Locating Where Faults Will Be," *Proc. Conf. Diversity in Computing*, pp. 48-50, 2005. (Paper=131, Status=F, Phase=2).
- [132] T.J. Ostrand, E.J. Weyuker, and R.M. Bell, "Automating Algorithms for the Identification of Fault-Prone Files," *Proc. Int'l Symp. Software Testing and Analysis*, pp. 219-227, 2007. (Paper=132, Status=F, Phase=4).
- [133] T.J. Ostrand, E.J. Weyuker, and R. Bell, "Where the Bugs Are," *ACM SIGSOFT Software Eng. Notes*, vol. 29, pp. 86-96, 2004. (Paper=133, Status=P).
- [134] T.J. Ostrand, E.J. Weyuker, and R. Bell, "Predicting the Location and Number of Faults in Large Software Systems," *IEEE Trans. Software Eng.*, vol. 31, no. 4, pp. 340-355, Apr. 2005. (Paper=134, Status=F, Phase=4).

- [135] T. Ostrand, E. Weyuker, and R. Bell, "Programmer-Based Fault Prediction," *Proc. Sixth Int'l Conf. Predictive Models in Software Eng.*, pp. 1-10, 2010. (Paper=135, Status=P).
- [136] T. Ostrand, E. Weyuker, R. Bell, and R. Ostrand, "A Different View of Fault Prediction," *Proc. 29th Ann. Int'l Computer Software and Applications Conf.*, vol. 2, pp. 3-4, July 2005. (Paper=136, Status=F, Phase=1).
- [137] F. Padberg, T. Ragg, and R. Schoknecht, "Using Machine Learning for Estimating the Defect Content After an Inspection," *IEEE Trans. Software Eng.*, vol. 30, no. 1, pp. 17-28, Jan. 2004. (Paper=137, Status=F, Phase=2).
- [138] G. Pai and J. Dugan, "Empirical Analysis of Software Fault Content and Fault Proneness Using Bayesian Methods," *IEEE Trans. Software Eng.*, vol. 33, no. 10, pp. 675-686, Oct. 2007. (Paper=138, Status=F, Phase=2, Data=N).
- [139] A.K. Pandey and N.K. Goyal, "Test Effort Optimization by Prediction and Ranking of Fault-Prone Software Modules," *Proc. Second Int'l Conf. Reliability, Safety and Hazard*, pp. 136-142, 2010. (Paper=139, Status=F, Phase=2, Data=N).
- [140] L. Pelayo and S. Dick, "Applying Novel Resampling Strategies to Software Defect Prediction," *Proc. Ann. Meeting of the North Amer. Fuzzy Information Processing Soc.*, pp. 69-72, June 2007. (Paper=140, Status=F, Phase=2, Data=N).
- [141] P. Pendharkar, "Exhaustive and Heuristic Search Approaches for Learning a Software Defect Prediction Model," *J. Eng. Applications of Artificial Intelligence*, vol. 23, no. 1, pp. 34-40, 2010. (Paper=141, Status=F, Phase=5, Data=N).
- [142] H. Peng and Z. Jie, "Predicting Defect-Prone Software Modules at Different Logical Levels," *Proc. Int'l Conf. Research Challenges in Computer Science*, pp. 37-40, 2009. (Paper=142, Status=F, Phase=2, Data=N).
- [143] M. Pighin and A. Marzona, "An Empirical Analysis of Fault Persistence through Software Releases," *Proc. Int'l Symp. Empirical Software Eng.*, pp. 206-212, Sept.-Oct. 2003. (Paper=143, Status=F, Phase=1).
- [144] M. Pinzger, N. Nagappan, and B. Murphy, "Can Developer-Module Networks Predict Failures?" *Proc. 16th ACM SIGSOFT Int'l Symp. Foundations of Software Eng.*, pp. 2-12, 2008. (Paper=144, Status=F, Phase=2).
- [145] R. Ramler, S. Lardorfer, and T. Natschlagler, "What Software Repositories Should Be Mined for Defect Predictors?" *Proc. 35th Euromicro Conf. Software Eng. and Advanced Applications*, pp. 181-187, 2009. (Paper=145, Status=F, Phase=4).
- [146] Z. Rana, S. Shamil, and M. Awais, "Ineffectiveness of Use of Software Science Metrics as Predictors of Defects in Object Oriented Software," *Proc. WRI World Congress Software Eng.*, vol. 4, pp. 3-7, May 2009. (Paper=146, Status=F, Phase=1, Data=N).
- [147] J. Ratzinger, T. Sigmund, and H.C. Gall, "On the Relation of Refactorings and Software Defect Prediction," *Proc. Int'l Working Conf. Mining Software Repositories*, pp. 35-38, 2008. (Paper=147, Status=F, Phase=4).
- [148] M. Reformat, *A Fuzzy-Based Meta-model for Reasoning about the Number of Software Defects*. Springer, 2003. (Paper=148, Status=F, Phase=2).
- [149] D. Rodriguez, R. Ruiz, J. Cuadrado-Gallego, and J. Aguilar-Ruiz, "Detecting Fault Modules Applying Feature Selection to Classifiers," *Proc. IEEE Int'l Conf. Information Reuse and Integration*, pp. 667-672, Aug. 2007. (Paper=149, Status=F, Phase=2, Data=N).
- [150] P.S. Sandhu, R. Goel, A.S. Brar, J. Kaur, and S. Anand, "A Model for Early Prediction of Faults in Software Systems," *Proc. Second Int'l Conf. Computer and Automation Eng.*, vol. 4, pp. 281-285, 2010. (Paper=150, Status=F, Phase=2, Data=N).
- [151] P.S. Sandhu, M. Kaur, and A. Kaur, "A Density Based Clustering Approach for Early Detection of Fault Prone Modules," *Proc. Int'l Conf. Electronics and Information Eng.*, vol. 2, pp. V2-525-V2-530, 2010. (Paper=151, Status=F, Phase=2, Data=N).
- [152] N. Schneidewind, "Investigation of Logistic Regression as a Discriminant of Software Quality," *Proc. Seventh Int'l Software Metrics Symp.*, pp. 328-337, 2001. (Paper=152, Status=F, Phase=2).
- [153] A. Schröter, T. Zimmermann, R. Premraj, and A. Zeller, "If Your Bug Database Could Talk," *Proc. Fifth Int'l Symp. Empirical Software Eng.*, vol. 2, pp. 18-20, 2006. (Paper=153, Status=F, Phase=1).
- [154] A. Schröter, T. Zimmermann, and A. Zeller, "Predicting Component Failures at Design Time," *Proc. ACM/IEEE Int'l Symp. Empirical Software Eng.*, pp. 18-27, 2006. (Paper=154, Status=P).
- [155] A. Schröter, T. Zimmermann, R. Premraj, and A. Zeller, "Where Do Bugs Come From?" *SIGSOFT Software Eng. Notes*, vol. 31, pp. 1-2, Nov. 2006. (Paper=155, Status=F, Phase=1).
- [156] C. Seiffert, T.M. Khoshgoftaar, and J.V. Hulse, "Improving Software-Quality Predictions with Data Sampling and Boosting," *IEEE Trans. Systems, Man, and Cybernetics, Part A: Systems and Humans*, vol. 39, no. 6, pp. 1283-1294, Nov. 2009. (Paper=156, Status=F, Phase=2, Data=N).
- [157] N. Seliya, T.M. Khoshgoftaar, and J. Van Hulse, "Predicting Faults in High Assurance Software," *Proc. IEEE 12th Int'l Symp. High-Assurance Systems Eng.*, pp. 26-34, 2010. (Paper=157, Status=F, Phase=2, Data=N).
- [158] N. Seliya, T. Khoshgoftaar, and S. Zhong, "Analyzing Software Quality with Limited Fault-Proneness Defect Data," *Proc. IEEE Ninth Int'l Symp. High-Assurance Systems Eng.*, pp. 89-98, Oct. 2005. (Paper=158, Status=F, Phase=2, Data=N).
- [159] R. Selvarani, T. Nair, and V. Prasad, "Estimation of Defect Proneness Using Design Complexity Measurements in Object-Oriented Software," *Proc. Int'l Conf. Signal Processing Systems*, pp. 766-770, May 2009. (Paper=159, Status=F, Phase=1).
- [160] R. Shatnawi and W. Li, "The Effectiveness of Software Metrics in Identifying Error-Prone Classes in Post-Release Software Evolution Process," *J. Systems and Software*, vol. 81, no. 11, pp. 1868-1882, 2008. (Paper=160, Status=P).
- [161] M. Sherriff, S.S. Heckman, M. Lake, and L. Williams, "Identifying Fault-Prone Files Using Static Analysis Alerts through Singular Value Decomposition," *Proc. Conf. Center for Advanced Studies on Collaborative Research*, pp. 276-279, 2007. (Paper=161, Status=F, Phase=2).
- [162] M. Sherriff, N. Nagappan, L. Williams, and M. Vouk, "Early Estimation of Defect Density Using an In-Process Haskell Metrics Model," *SIGSOFT Software Eng. Notes*, vol. 30, pp. 1-6, May 2005. (Paper=162, Status=F, Phase=1).
- [163] Y. Shin, R.M. Bell, T.J. Ostrand, and E.J. Weyuker, "Does Calling Structure Information Improve the Accuracy of Fault Prediction?" *Proc. Sixth Int'l Working Conf. Mining Software Repositories*, pp. 61-70, 2009. (Paper=163, Status=P).
- [164] S. Shivaji, E.J. Whitehead, R. Akella, and K. Sunghun, "Reducing Features to Improve Bug Prediction," *Proc. IEEE/ACM 24th Int'l Conf. Automated Software Eng.*, pp. 600-604, 2009. (Paper=164, Status=P).
- [165] P. Singh and S. Verma, "An Investigation of the Effect of Discretization on Defect Prediction Using Static Measures," *Proc. Int'l Conf. Advances in Computing, Control, Telecomm. Technologies*, pp. 837-839, 2009. (Paper=165, Status=F, Phase=2).
- [166] Y. Singh, A. Kaur, and R. Malhotra, "Predicting Software Fault Proneness Model Using Neural Network," *Product-Focused Software Process Improvement*, vol. 5089, pp. 204-214, 2008. (Paper=166, Status=F, Phase=2, Data=N).
- [167] Y. Singh, A. Kaur, and R. Malhotra, "Empirical Validation of Object-Oriented Metrics for Predicting Fault Proneness Models," *Software Quality J.*, vol. 18, no. 1, pp. 3-35, 2010. (Paper=167, Status=F, Phase=2, Data=N).
- [168] Q. Song, M. Shepperd, M. Cartwright, and C. Mair, "Software Defect Association Mining and Defect Correction Effort Prediction," *IEEE Trans. Software Eng.*, vol. 32, no. 2, pp. 69-82, Feb. 2006. (Paper=168, Status=F, Phase=2).
- [169] C. Stringfellow and A. Andrews, "Deriving a Fault Architecture to Guide Testing," *Software Quality J.*, vol. 10, no. 4, pp. 299-330, 2002. (Paper=169, Status=F, Phase=1).
- [170] G. Succi, W. Pedrycz, M. Stefanovic, and J. Miller, "Practical Assessment of the Models for Identification of Defect-Prone Classes in Object-Oriented Commercial Systems Using Design Metrics," *J. Systems and Software*, vol. 65, no. 1, pp. 1-12, 2003. (Paper=170, Status=F, Phase=1).
- [171] M.D.M. Suffian and M.R. Abdullah, "Establishing a Defect Prediction Model Using a Combination of Product Metrics as Predictors via Six Sigma Methodology," *Proc. Int'l Symp. Information Technology*, vol. 3, pp. 1087-1092, 2010. (Paper=171, Status=F, Phase=2).
- [172] M.M.T. Thwin and T.-S. Quah, "Application of Neural Network for Predicting Software Development Faults Using Object-Oriented Design Metrics," *Proc. Ninth Int'l Conf. Neural Information Processing*, vol. 5, pp. 2312-2316, Nov. 2002. (Paper=172, Status=F, Phase=2).

- [173] P. Tomaszewski, H. Grahm, and L. Lundberg, "A Method for an Accurate Early Prediction of Faults in Modified Classes," *Proc. IEEE 22nd Int'l Conf. Software Maintenance*, pp. 487-496, Sept. 2006. (Paper=173, Status=F, Phase=2).
- [174] A. Tosun and A. Bener, "Reducing False Alarms in Software Defect Prediction by Decision Threshold Optimization," *Proc. Third Int'l Symp. Empirical Software Eng. and Measurement*, pp. 477-480, 2009. (Paper=174, Status=F, Phase=2, Data=N).
- [175] A. Tosun, B. Turhan, and A. Bener, "Practical Considerations in Deploying AI for Defect Prediction: A Case Study within the Turkish Telecommunication Industry," *Proc. Fifth Int'l Conf. Predictor Models in Software Eng.*, p. 11, 2009. (Paper=175, Status=F, Phase=2, Data=N).
- [176] A. Tosun, A.B. Bener, B. Turhan, and T. Menzies, "Practical Considerations in Deploying Statistical Methods for Defect Prediction: A Case Study within the Turkish Telecommunications Industry," *Information and Software Technology*, vol. 52, no. 11, pp. 1242-1257, 2010. (Paper=176, Status=F, Phase=2, Data=N).
- [177] B. Turhan and A. Bener, "A Multivariate Analysis of Static Code Attributes for Defect Prediction," *Proc. Seventh Int'l Conf. Quality Software*, pp. 231-237, Oct. 2007. (Paper=177, Status=F, Phase=2, Data=N).
- [178] B. Turhan, G. Kocak, and A. Bener, "Software Defect Prediction Using Call Graph Based Ranking (CGBR) Framework," *Proc. 34th Euromicro Conf. Software Eng. and Advanced Applications*, pp. 191-198, Sept. 2008. (Paper=178, Status=F, Phase=2).
- [179] B. Turhan, G. Kocak, and A. Bener, "Data Mining Source Code for Locating Software Bugs: A Case Study in Telecommunication Industry," *Expert Systems with Applications*, vol. 36, no. 6, pp. 9986-9990, 2009. (Paper=179, Status=F, Phase=2, Data=N).
- [180] B. Turhan, A.B. Bener, and T. Menzies, "Regularities in Learning Defect Predictors," *Proc. 11th Int'l Conf. Product-Focused Software Process Improvement*, pp. 116-130, 2010. (Paper=180, Status=F, Phase=4, Data=N).
- [181] B. Turhan, T. Menzies, A. Bener, and J. Di Stefano, "On the Relative Value of Cross-Company and within-Company Data for Defect Prediction," *Empirical Software Eng.*, vol. 14, no. 5, pp. 540-578, 2009. (Paper=181, Status=F, Phase=2, Data=N).
- [182] O. Vandecruys, D. Martens, B. Baesens, C. Mues, M. De Backer, and R. Haesen, "Mining Software Repositories for Comprehensible Software Fault Prediction Models," *J. Systems and Software*, vol. 81, no. 5, pp. 823-839, 2008. (Paper=182, Status=F, Phase=2, Data=N).
- [183] R. Vivanco, Y. Kamei, A. Monden, K. Matsumoto, and D. Jin, "Using Search-Based Metric Selection and Oversampling to Predict Fault Prone Modules," *Proc. 23rd Canadian Conf. Electrical and Computer Eng.*, pp. 1-6, 2010. (Paper=183, Status=F, Phase=2, Data=N).
- [184] D. Wahyudin, A. Schatten, D. Winkler, A.M. Tjoa, and S. Biffl, "Defect Prediction Using Combined Product and Project Metrics—A Case Study from the Open Source 'Apache' Myfaces Project Family," *Proc. 34th Euromicro Conf. Software Eng. and Advanced Applications*, pp. 207-215, 2008. (Paper=184, Status=F, Phase=1).
- [185] T. Wang and W.-h. Li, "Naive Bayes Software Defect Prediction Model," *Proc. Int'l Conf. Computational Intelligence and Software Eng.*, pp. 1-4, 2010. (Paper=185, Status=F, Phase=2, Data=N).
- [186] W. Wei, D. Xuan, L. Chunping, and W. Hui, "A Novel Evaluation Method for Defect Prediction in Software Systems," *Proc. Int'l Conf. Computational Intelligence and Software Eng.*, pp. 1-5, 2010. (Paper=186, Status=F, Phase=1).
- [187] Y. Weimin and L. Longshu, "A Rough Set Model for Software Defect Prediction," *Proc. Int'l Conf. Intelligent Computation Technology and Automation*, vol. 1, pp. 747-751, 2008. (Paper=187, Status=F, Phase=2, Data=N).
- [188] E. Weyuker, T. Ostrand, and R. Bell, "Using Developer Information as a Factor for Fault Prediction," *Proc. Third Int'l Workshop Predictor Models in Software Eng.*, p. 8, May 2007. (Paper=188, Status=F, Phase=2).
- [189] E. Weyuker, T. Ostrand, and R. Bell, "Comparing Negative Binomial and Recursive Partitioning Models for Fault Prediction," *Proc. Fourth Int'l Workshop Predictor Models in Software Eng.*, pp. 3-10, 2008. (Paper=189, Status=F, Phase=2).
- [190] E. Weyuker, T. Ostrand, and R. Bell, "Do Roo Many Cooks Spoil the Broth? Using the Number of Developers to Enhance Defect Prediction Models," *Empirical Software Eng.*, vol. 13, no. 5, pp. 539-559, 2008. (Paper=190, Status=P).
- [191] E.J. Weyuker, T.J. Ostrand, and R.M. Bell, "Comparing the Effectiveness of Several Modeling Methods for Fault Prediction," *Empirical Software Eng.*, vol. 15, no. 3, pp. 277-295, 2010. (Paper=191, Status=F, Phase=2).
- [192] C. Wohlin, M. Host, and M. Ohlsson, "Understanding the Sources of Software Defects: A Filtering Approach," *Proc. Eighth Int'l Workshop Program Comprehension*, pp. 9-17, 2000. (Paper=192, Status=F, Phase=1).
- [193] W. Wong, J. Horgan, M. Syring, W. Zage, and D. Zage, "Applying Design Metrics to Predict Fault-Prone: A Case Study on a Large-Scale Software System," *Software: Practice and Experience*, vol. 30, no. 14, pp. 1587-1608, 2000. (Paper=193, Status=F, Phase=2).
- [194] Z. Xu, T. Khoshgoftaar, and E. Allen, "Prediction of Software Faults Using Fuzzy Nonlinear Regression Modeling," *Proc. IEEE Fifth Int'l Symp. High Assurance Systems Eng.*, pp. 281-290, 2000. (Paper=194, Status=F, Phase=2).
- [195] B. Yang, L. Yao, and H.-Z. Huang, "Early Software Quality Prediction Based on a Fuzzy Neural Network Model," *Proc. Third Int'l Conf. Natural Computation*, vol. 1, pp. 760-764, Aug. 2007. (Paper=195, Status=F, Phase=2).
- [196] L. Yi, T.M. Khoshgoftaar, and N. Seliya, "Evolutionary Optimization of Software Quality Modeling with Multiple Repositories," *IEEE Trans. Software Eng.*, vol. 36, no. 6, pp. 852-864, Nov. 2010. (Paper=196, Status=F, Phase=2, Data=N).
- [197] H. Youngki, B. Jongmoon, K. In-Young, and C. Ho-Jin, "A Value-Added Predictive Defect Type Distribution Model Based on Project Characteristics," *Proc. IEEE/ACIS Seventh Int'l Conf. Computer and Information Science*, pp. 469-474, 2008. (Paper=197, Status=F, Phase=2).
- [198] P. Yu, T. Systa, and H. Muller, "Predicting Fault-Prone: Using OO Metrics an Industrial Case Study," *Proc. Sixth European Conf. Software Maintenance and Reeng.*, pp. 99-107, 2002. (Paper=198, Status=F, Phase=1).
- [199] X. Yuan, T. Khoshgoftaar, E. Allen, and K. Ganesan, "An Application of Fuzzy Clustering to Software Quality Prediction," *Proc. IEEE Third Symp. Application-Specific Systems and Software Eng. Technology*, pp. 85-90, 2000. (Paper=199, Status=F, Phase=2).
- [200] H. Zhang, A. Nelson, and T. Menzies, "On the Value of Learning from Defect Dense Components for Software Defect Prediction," *Proc. Sixth Int'l Conf. Predictive Models in Software Eng.*, p. 14, Sept. 2010. (Paper=200, Status=F, Phase=2, Data=N).
- [201] S. Zhong, T. Khoshgoftaar, and N. Seliya, "Unsupervised Learning for Expert-Based Software Quality Estimation," *Proc. IEEE Eighth Int'l Symp. High Assurance Systems Eng.*, pp. 149-155, Mar. 2004. (Paper=201, Status=F, Phase=2, Data=N).
- [202] Y. Zhou and H. Leung, "Empirical Analysis of Object-Oriented Design Metrics for Predicting High and Low Severity Faults," *IEEE Trans. Software Eng.*, vol. 32, no. 10, pp. 771-789, Oct. 2006. (Paper=202, Status=F, Phase=2, Data=N).
- [203] Y. Zhou, B. Xu, and H. Leung, "On the Ability of Complexity Metrics to Predict Fault-Prone Classes in Object-Oriented Systems," *J. Systems and Software*, vol. 83, no. 4, pp. 660-674, 2010. (Paper=203, Status=P).
- [204] T. Zimmermann and N. Nagappan, "Predicting Subsystem Failures Using Dependency Graph Complexities," *Proc. IEEE 18th Int'l Symp. Software Reliability*, pp. 227-236, Nov. 2007. (Paper=204, Status=F, Phase=2).
- [205] T. Zimmermann and N. Nagappan, "Predicting Defects Using Network Analysis on Dependency Graphs," *Proc. ACM/IEEE 30th Int'l Conf. Software Eng.*, pp. 531-540, 2008. (Paper=205, Status=F, Phase=4).
- [206] T. Zimmermann and N. Nagappan, "Predicting Defects with Program Dependencies," *Proc. Third Int'l Symp. Empirical Software Eng. and Measurement*, pp. 435-438, 2009. (Paper=206, Status=F, Phase=2).
- [207] T. Zimmermann, R. Premraj, and A. Zeller, "Predicting Defects for Eclipse," *Proc. Int'l Workshop Predictor Models in Software Eng.*, p. 9, May 2007. (Paper=207, Status=F, Phase=2).
- [208] T. Zimmermann, N. Nagappan, H. Gall, E. Giger, and B. Murphy, "Cross-Project Defect Prediction: A Large Scale Experiment on Data vs. Domain vs. Process," *Proc. Seventh Joint Meeting of the European Software Eng. Conf. and the ACM SIGSOFT Symp. The Foundations of Software Eng.*, pp. 91-100, 2009. (Paper=208, Status=F, Phase=4).

TABLE 10
Conferences and Journals Manually Searched

Conference manually searched	Journals manually searched
International Conference on Software Engineering (ICSE)	IEEE Transactions of Software Engineering
International Conference on Software Maintenance (ICSM)	Journal of Systems and Software
IEEE Int'l Working Conference on Source Code Analysis and Manipulation (SCAM)	Journal of Empirical Software Engineering
International Conference on Automated Software Engineering	Software Quality Journal
IEEE Int'l Symposium and Workshop on Engineering of Computer Based Systems	Information & Software Technology
International Symposium on Automated Analysis-driven Debugging	
International Symposium on Software Testing and Analysis (ISSTA)	
International Symposium on Software Reliability Engineering	
ACM SIGPLAN Conference on Programming language Design and Implementation	
Int'l Workshop on Mining Software Repositories	
Empirical Software Engineering & Measurement	
PROMISE	
Foundations of Software Engineering	

TABLE 11
Additional Data Quality Criteria

Data quality criteria	Criteria definitions	Why the criteria is important
Has any data cleaning been done?	An indication that the quality of the data has been considered is necessary and that any data cleaning needed has been addressed e.g. missing values handled, outliers and errorful data been removed.	Data sets are often noisy. They often contain outliers and missing values that can skew results (the impact of this depends on the analysis methods used). Our confidence in the predictions made by a model is impacted by the quality of the data used while building the model. Few studies have cleaned their data and so we did not apply this criterion.
Have repeated attributes been removed?	An indication that the impact of repeated attributes has been considered should be given. For example machine learning studies could mention attribute selection while other studies could consider, for example, Principal Component Analysis.	Repeated attributes and related attributes have been shown to bias the outcomes of models. Confidence is affected in the predictions of studies which have not considered the impact of repeated/related attributes. Few studies have considered repeated attributes and so we did not apply this criterion.

APPENDIX A

SEARCH STRING

The following search string was used in our searches:

(Fault* OR bug* OR defect* OR errors OR corrections OR corrective OR fix*) *in title only*

AND (Software) *anywhere in study*

APPENDIX B

CONFERENCES AND JOURNALS MANUALLY SEARCHED

See Table 10

APPENDIX C

ADDITIONAL ASSESSMENT CRITERIA

Data quality criteria. The efficacy of the predictions made by a model is determined by the quality of the data on which the model was built. Leibchen and Shepperd [26] report that many studies do not seem to consider the quality of the data they use. Many fault prediction models are based on machine learning, where it has been shown that a lack of data cleaning may compromise the predictions obtained [21]. The criteria shown in Table 11 are based on [21], [S168], [S192], [S194], and [S19].

Predictive performance criteria. Measuring the predictive performance of a model is an essential part of demonstrating the usefulness of that model. Measuring model performance is complex and there are many ways in

which the performance of a model may be measured. Furthermore, the value of measures varies according to context. For example, safety critical system developers may want models that identify as many faults as possible, accepting the cost of false alarms, whereas business system developers may want models which do not generate many false alarms as testing effort is short to ensure the timely release of a product at the cost of missing some faults. Appendix D reports the principles of predictive performance measurement and provides the basis of our performance measurement criteria. Table 12 shows our predictive performance measurement criteria.

APPENDIX D

THE PRINCIPLES OF PREDICTIVE PERFORMANCE MEASUREMENT

This overview of measuring predictive performance is based on [30], [S61], and [S97]. The measurement of predictive performance is often based on the analysis of data in a confusion matrix (shown in Table 13 and explained further in Table 14). This matrix reports how the model classified the different fault categories compared to their actual classification (predicted versus observed). Many performance measures are related to components of the confusion matrix shown in Table 14. Confusion matrix-based measures are most relevant to fault prediction models producing categorical outputs, though continuous outputs can be converted to categorical outputs and analyzed in terms of a confusion matrix.

TABLE 12
Additional Predictive Performance Measurement Criteria

Measurement criteria	Criteria definitions	Why the criteria is important
Have imbalanced data sets been accounted for, sufficient to enable confidence in predictive performance?	No one approach adequately accounts for imbalanced data in every circumstance. See Appendix F for an overview of the approaches available. Each study should be assessed on a case-by-case basis according to the specific model reported.	See Appendix F for an overview of the importance of dealing with imbalanced data. There is no one best way of accounting for data imbalance and it was thus difficult to identify a precise enough criterion to apply consistently across studies. In addition there remains significant debate on data imbalance (see [27], [[179]], [28], [29]).
Has predictive performance been reported appropriately?	For each model reporting categorical results a study should report either: - A confusion matrix - Area Under the Curve (AUC) For each model reporting continuous results a study choosing to report measures of error should not report only Mean Squared Error. Average Relative Error should also be reported. Or else results based on Chi Square should be reported.	The particular set of performance measures reported by studies can make it difficult to understand how a model performs overall in terms of correct and incorrect predictions. Confusion matrix constructs form the basis of most other ways of reporting predictive performance. Reporting the confusion matrix (possibly in addition to other measures reported by studies) would allow subsequent analysis of performance in ways other than those preferred by the model developer. Menzies et al. [[114]] suggests a useful way in which data from multiple confusion matrices may be effectively reported. AUC performance data is reported to be an effective way in which to compare the ability of a modelling technique to cope with different datasets (Lessmann et al. [[97]]). Reporting AUC means that models using different datasets can then be meta-analysed. This would enable a much richer understanding of the abilities of particular modelling techniques. Ideally data for the whole ROC curve would be given by each study. It is impractical to report this amount of data and would require the use of on-line data stores. However AUC has limitations for imbalanced data sets as reported by [23] [12]. Mean Squared Error (MSE) generates results that can only be interpreted within the data set from which they originated as the measurement scales used dictate the size of the error. MSE limits the comparability of results across other data sets. Chi Square or Average Relative Error should be used instead. Only a small number of models currently report confusion matrix and AUC data. Consequently, applying this criterion was untenable as so few studies would have passed. Similarly MSE is a poorly understood measure and its limitations are not widely reported and so the current application of this requirement is not tenable.

Composite performance measures can be calculated by combining values from the confusion matrix (see Table 15). “Recall” (otherwise known as the true positive rate, probability of detection (pd), or sensitivity) describes the proportion of faulty code units (usually files, modules, or packages) correctly predicted as such, while “precision” describes how reliable a prediction is in terms of what proportion of code predicted as faulty actually was faulty. Both are important when test sets are imbalanced, but there is a tradeoff between these two measures [S61]. An additional composite measure is the false positive rate (pf), which describes the proportion of erroneous defective predictions. Thus, the optimal classifier would achieve a pd of 1, precision of 1, and a pf of 0. The performance measure balance combines pd and pf. A high-balance value (near 1) is achieved with a high pd and low pf. Balance can also be adjusted to a factor in the cost of false alarms which typically do not result in fault fixes. When the combinations of pd and pf are plotted, they produce a Receiver Operator

Curve (ROC). This gives a range of balance figures, and it is usual to report the area under the curve as varying between 0 and 1, with 1 being the ideal value. Table 16 shows other ways in which the performance of a model can be measured. Such measures are usually used in models that produce continuous or ranking results.

APPENDIX E

CALCULATING PRECISION, RECALL, AND F-MEASURE FOR CATEGORICAL STUDIES (REPORTED IN [31])

Many studies report precision and recall, but others report pd and pf. If we are to compare the results we need to convert the results of one paper into the performance

TABLE 13
Confusion Matrix

	Predicted defective	Predicted defect free
Observed defective	True Positive (TP)	False Negative (FN)
Observed defect free	False Positive (FP)	True Negative (TN)

TABLE 14
Confusion Matrix-Based Performance Indicator

Construct	Also known as	Description
False Positive	FP, and Type I Error	Classifies non faulty unit as faulty
False Negative	FN, and Type II Error	Classifies faulty unit as not faulty
True Positive	TP	Correctly classified as faulty
True Negative	TN	Correctly classified as non-faulty

TABLE 15
Composite Performance Measures

Construct	Defined as	Description
Recall pd (probability of detection) Sensitivity True positive rate	$TP/(TP + FN)$	Proportion of faulty units correctly classified
Precision pf (probability of false alarm) False positive rate	$TP/(TP + FP)$ $FP/(FP + TN)$	Proportion of units correctly predicted as faulty Proportion of non-faulty units incorrectly classified
Specificity True negative rate	$TN/(TN + FP)$	Proportion of correctly classified non faulty units
f-measure	$\frac{2 \cdot Recall \cdot Precision}{Recall + Precision}$ $(TN + TP)$	Most commonly defined as the harmonic mean of precision and recall
Accuracy	$\frac{TP + TN}{(TN + FN + FP + TP)}$	Proportion of correctly classified units
Mis-classification rate Error-rate	$1 - accuracy$	Proportion of incorrectly classified units
Balance	$1 - \frac{\sqrt{(0 - pf)^2 + (1 - pd)^2}}{\sqrt{2}}$	Combines <i>pf</i> and <i>pd</i> into one measure and is most commonly defined as the distance from the ROC 'sweet spot' (where <i>pd</i> =1, <i>pf</i> =0).
Receiver operating characteristic (ROC) curve		A graphical plot of the sensitivity (or <i>pd</i>) vs. 1 - specificity (or <i>pf</i>) for a binary classification system where its discrimination threshold is varied

TABLE 16
Performance Indicators Defined

Measure	Constructs and Definitions
Error measures	Average residual error, relative error, relative square error, standard error of estimate, root mean squared error, median relative error, mean square error, mean absolute error, mean absolute relative error, error rate.
Significance of difference between predicted and observed	Spearmans, Pearsons, Chi Square

measures reported by the other paper. In this case, we want to report everything in terms of precision and recall. We chose these measures as fault prediction datasets are often highly imbalanced (Zhang and Zhang [27] and Gray et al. [12]). When trying to compare the results of one paper with the results of another paper, it may be necessary to reconstruct a form of the Confusion Matrix (see Table 13 in Appendix D) where the values are not the sums of instances, but the frequency of each instance:

$$1 = TP + TN + FP + FN. \quad (1)$$

This is possible in many cases when the distribution of the classes is also reported. To do this we need to know the frequency of the true class *d*, where

$$d = TP + FN. \quad (2)$$

It then becomes possible to calculate *TP*, *FP*, *TN*, and *FN* as follows:

Given *pf* and *d*:

$$TN = (1 - d)(1 - pf), \quad (3)$$

$$FP = (1 - d)pf. \quad (4)$$

Given *pd*(*Recall*(*r*)) and *d*:

$$TP = d \cdot r, \quad (5)$$

$$FN = d(1 - r). \quad (6)$$

Given *FNR*(*TypeII*(*t2*)), *pf* and *d* we already have (1), (3), and (4):

$$FN = \frac{pf(1 - d)t2}{(1 - t2)}, \quad (7)$$

$$TP = 1 - FN - TN - FP. \quad (8)$$

Given *Precision*(*p*), *Recall*(*r*), and *d* we already have (1), (5), and (6):

$$FP = \frac{FN(1 - p)}{p} = \frac{d(1 - r)(1 - p)}{p}, \quad (9)$$

$$TN = 1 - FP - FN - TP. \quad (10)$$

In some cases *d* is not available but more performance measures are provided.

Given *Errorrate*(*er*), *FNR*(*TypeII*(*t2*)), and *pf*:

$$d = 1 - \frac{er(1 - t2)}{pf}, \quad (11)$$

which can then be used with (3), (4), (7), and (8).

Given *Precision*(*p*), *Recall*(*r*), and *Accuracy*(*a*):

$$d = \frac{p(1 - a)}{p - 2pr + r}, \quad (12)$$

which can then be used with (5), (6), (9), and (10).

Given $Accuracy(a)$, pf , and $FNR(TypeII(t2))$:

$$FP = \frac{(1 - t2 - a)pf}{pf - t2}, \quad (13)$$

$$TN = \frac{(1 - t2 - a)(1 - pf)}{pf - t2}, \quad (14)$$

$$TP = \frac{(1 - t2)(pf - 1 + 1)}{pf - t2}, \quad (15)$$

$$FN = 1 - TP - TN - FP. \quad (16)$$

Given FP , FN , and d :

$$TP = d - FP, \quad (17)$$

which can then be used with (10).

The following values were extracted from [S83]:

$$er = 0.3127, pf = 0.3134, t2 = 0.2826.$$

We compute

$$d = 0.2842.$$

Giving,

$$FN = 0.0884, TN = 0.4915, FP = 0.2243, TP = 0.1958.$$

Finally,

$$\begin{aligned} Precision &= 0.4661, Recall = 0.6891, \\ F\text{-measure} &= 0.5561. \end{aligned}$$

APPENDIX F

THE CLASS IMBALANCE PROBLEM

Substantially imbalanced datasets are commonly used in binary fault prediction studies (i.e., there are usually many more nonfaulty units than faulty units) [32], [27]. An extreme example of this is seen in NASA dataset PC2, which has only 0.4 percent of data points belonging to the faulty class (23 out of 5,589 data points). This distribution of faulty and nonfaulty units—known as the class distribution—should be taken into account during any binary fault prediction task. This is because imbalanced data can strongly influence both the training of a classification model and the suitability of classifier performance metrics.

When training a classifier using imbalanced data, an algorithm can struggle to learn from the minority class. This is typically due to an insufficient quantity of minority class data. The most common symptom when this occurs is for a classifier to predict all data points as belonging to the majority class, which is of little practical worth. To avoid this happening, various approaches can be used and are typically based around training-set sampling and/or learning algorithm optimization. Note that these techniques are entirely optional, and may not be necessary. This is because learning techniques vary in their sensitivity to imbalanced data. For example, C4.5 decision trees have been reported to struggle with imbalanced data [16] and [17], whereas fuzzy-based classifiers have been reported to perform robustly regardless of class distribution [33].

Sampling methods involve the manipulation of training data in order to reduce the level of imbalance and therefore alleviate the problems associated with learning from imbalanced data. Undersampling methods involve reducing the size of the majority class, whereas oversampling methods involve increasing the size of the minority class. Such techniques have been reported to be useful [11]; however, they do suffer from drawbacks. With undersampling methods, the main problem is deciding which majority class data points should be removed. With oversampling methods, there is a risk of the learning algorithm overfitting the oversampled data. This will probably result in good training data performance, but low performance when the classifier is presented with unseen data (data independent from that used during training) [11].

Many learning algorithms can have their various parameters adjusted in order to boost performance on imbalanced data. This can be very effective, as many algorithms by default assume an equal class distribution during training. By increasing the misclassification cost of the minority class, it is possible to construct models that are better suited to imbalanced domains. Such methods can be difficult and/or time consuming to approximate appropriate misclassification costs.

Additional problems caused by imbalanced data are that selecting appropriate classifier performance measures is more difficult. This is because measures which favor the majority class (such as accuracy and error rate) are no longer sufficient [11]. More appropriate measures in imbalanced domains include: precision, recall, f-measure (see Appendix D), and g-mean [11].

In contrast to the training data, the balance of test data should be representative of that which will be encountered in the real world.

There remains significant debate on data imbalance in fault prediction (see [12], [27], [S179], [28], [29]).

ACKNOWLEDGMENTS

The authors are grateful to the United Kingdom's Engineering and Physical Science Research Council who supported this research at Brunel University under grant EPSRC EP/E063039/1 and to Science Foundation Ireland grant 3/CE2/I303_1 who partially supported this work at Lero. They are also grateful to Dr. Sue Black and Dr. Paul Wernick who provided input to the early stages of the work reported in this paper and to Professor Martin Shepperd for his suggestions throughout the work. They are also grateful for the detailed and insightful comments from the reviewers that enabled them to significantly improve the quality of this paper.

REFERENCES

- [1] N. Fenton and M. Neil, "A Critique of Software Defect Prediction Models," *IEEE Trans. Software Eng.*, vol. 25, no. 5, pp. 675-689, Sept. 1999.
- [2] C. Catal and B. Diri, "A Systematic Review of Software Fault Prediction Studies," *Expert Systems with Application*, vol. 36, no. 4, pp. 7346-7354, 2009.
- [3] B. Kitchenham and S. Charters, "Guidelines for Performing Systematic Literature Reviews in Software Engineering (Version 2.3)," Technical Report EBSE-2007-01, Keele Univ., EBSE, 2007.

- [4] M. Jorgensen and M. Shepperd, "A Systematic Review of Software Development Cost Estimation Studies," *IEEE Trans. Software Eng.*, vol. 33, no. 1, pp. 33-53, Jan. 2007.
- [5] B. Kitchenham, "What's Up with Software Metrics?—A Preliminary Mapping Study," *J. Systems and Software*, vol. 83, no. 1, pp. 37-51, 2010.
- [6] Q. Song, Z. Jia, M. Shepperd, S. Ying, and J. Liu, "A General Software Defect-Proneness Prediction Framework," *IEEE Trans. Software Eng.*, vol. 37, no. 3, pp. 356-370, May 2011.
- [7] K. Petersen and C. Wohlin, "Context in Industrial Software Engineering Research," *Proc. Third Int'l Symp. Empirical Software Eng. and Measurement*, pp. 401-404, 2009.
- [8] P.G. Armour, "Beware of Counting LOC," *Comm. ACM*, vol. 47, pp. 21-24, Mar. 2004.
- [9] D. Bowes and T. Hall, "SLuRp: A Web Enabled Database for Effective Management of Systematic Literature Reviews," Technical Report 510, Univ. of Hertfordshire, 2011.
- [10] T. Menzies, A. Dekhtyar, J. Distefano, and J. Greenwald, "Problems with Precision: A Response to 'Comments on 'Data Mining Static Code Attributes to Learn Defect Predictors,'" *IEEE Trans. Software Eng.*, vol. 33, no. 9, pp. 637-640, Sept. 2007.
- [11] H. He and E. Garcia, "Learning from Imbalanced Data," *IEEE Trans. Knowledge and Data Eng.*, vol. 21, no. 9, pp. 1263-1284, Sept. 2008.
- [12] D. Gray, D. Bowes, N. Davey, Y. Sun, and B. Christianson, "Further Thoughts on Precision," *Proc. Evaluation and Assessment in Software Eng.*, 2011.
- [13] D.S. Cruzes and T. Dybå, "Research Synthesis in Software Engineering: A Tertiary Study," *Information Software Technology*, vol. 53, pp. 440-455, May 2011.
- [14] R. Rosenthal and M. DiMatteo, "Meta-Analysis: Recent Developments in Quantitative Methods for Literature Reviews," *Ann. Rev. Psychology*, vol. 52, no. 1, pp. 59-82, 2001.
- [15] B. Turhan, T. Menzies, A. Bener, and J. Di Stefano, "On the Relative Value of Cross-Company and Within-Company Data for Defect Prediction," *Empirical Software Eng.*, vol. 14, no. 5, pp. 540-578, 2009.
- [16] N. Japkowicz and S. Stephen, "The Class Imbalance Problem: A Systematic Study," *Intelligent Data Analysis*, vol. 6, no. 5, pp. 429-449, 2002.
- [17] W. Liu, S. Chawla, D.A. Cieslak, and N.V. Chawla, "A Robust Decision Tree Algorithm for Imbalanced Data Sets," *Proc. 10th SIAM Int'l Conf. Data Mining*, pp. 766-777, 2010.
- [18] C. Hsu, C. Chang, and C. Lin, "A Practical Guide to Support Vector Classification," technical report, Dept. of Computer Science and Information Eng., Nat'l Taiwan Univ., 2003.
- [19] N. Cristianini and J. Shawe-Taylor, *An Introduction to Support Vector Machines: And Other Kernel-Based Learning Methods*. Cambridge Univ. Press, 2006.
- [20] T. Hall, D. Bowes, G. Liebchen, and P. Wernick, "Evaluating Three Approaches to Extracting Fault Data from Software Change Repositories," *Proc. 11th Int'l Conf. Product-Focused Software Process Improvement*, pp. 107-115, 2010.
- [21] D. Gray, D. Bowes, N. Davey, Y. Sun, and B. Christianson, "The Misuse of the NASA Metrics Data Program Data Sets for Automated Software Defect Prediction," *Proc. Evaluation and Assessment in Software Eng.*, 2011.
- [22] N. Pizzi, A. Summers, and W. Pedrycz, "Software Quality Prediction Using Median-Adjusted Class Labels," *Proc. Int'l Joint Conf. Neural Networks*, vol. 3, pp. 2405-2409, 2002.
- [23] J. Davis and M. Goadrich, "The Relationship between Precision-Recall and ROC Curves," *Proc. 23rd Int'l Conf. Machine Learning*, pp. 233-240, 2006.
- [24] A. Lamkanfi, S. Demeyer, E. Giger, and B. Goethals, "Predicting the Severity of a Reported Bug," *Proc. IEEE Working Conf. Seventh Mining Software Repositories*, pp. 1-10, 2010.
- [25] I. Myrtveit, E. Stensrud, and M. Shepperd, "Reliability and Validity in Comparative Studies of Software Prediction Models," *IEEE Trans. Software Eng.*, vol. 31, no. 5, pp. 380-391, May 2005.
- [26] G. Liebchen and M. Shepperd, "Data Sets and Data Quality in Software Engineering," *Proc. Fourth Int'l Workshop Predictor Models in Software Eng.*, pp. 39-44, 2008.
- [27] H. Zhang and X. Zhang, "Comments on 'Data Mining Static Code Attributes to Learn Defect Predictors,'" *IEEE Trans. Software Eng.*, vol. 33, no. 9, pp. 635-637, Sept. 2007.
- [28] G. Batista, R. Prati, and M. Monard, "A Study of the Behavior of Several Methods for Balancing Machine Learning Training Data," *ACM SIGKDD Explorations Newsletter*, vol. 6, no. 1, pp. 20-29, 2004.
- [29] Y. Kamei, A. Monden, S. Matsumoto, T. Kakimoto, and K. Matsumoto, "The Effects of Over and Under Sampling on Fault-Prone Module Detection," *Proc. First Int'l Symp. Empirical Software Eng. and Measurement*, pp. 196-204, Sept. 2007.
- [30] T. Ostrand and E. Weyuker, "How to Measure Success of Fault Prediction Models," *Proc. Fourth Int'l Workshop Software Quality Assurance: In Conjunction with the Sixth ESEC/FSE Joint Meeting*, pp. 25-30, 2007.
- [31] D. Bowes and D. Gray, "Recomputing the Confusion Matrix for Prediction Studies Reporting Categorical Output," Technical Report 509, Univ. of Hertfordshire, 2011.
- [32] N.V. Chawla, N. Japkowicz, and A. Kotcz, "Editorial: Special Issue on Learning from Imbalanced Data Sets," *SIGKDD Explorations*, vol. 6, no. 1, pp. 1-6, 2004.
- [33] S. Visa and A. Ralescu, "Fuzzy Classifiers for Imbalanced, Complex Classes of Varying Size," *Proc. Information Processing and Management of Uncertainty in Knowledge-Based Systems*, pp. 393-400, 2004.



ingly centered on fault prediction.



Tracy Hall received the PhD degree in software metrics from City University in 1998. Currently she is working as a reader in software engineering at Brunel University. Previously she was the head of the Systems & Software Research Group at the University of Hertfordshire. Over the last 15 years she has conducted many empirical software engineering studies with a variety of industrial collaborators. Her research interests have become increasingly centered on fault prediction.



Sarah Beecham is a research fellow working at Lero, The Irish Software Engineering Research Centre. She is currently collaborating closely with industry to develop a process model of global software development recommended practices. Her research interests include software engineer motivation, requirements engineering, fault prediction, effort estimation, and open source software development.

David Bowes is currently working toward the PhD degree, focusing on defect prediction studies. He has been a senior lecturer at the University of Hertfordshire since 2005. He has published mainly in the area of program slicing metrics, defect prediction, and epigenetic robotics.



David Gray received the BSc and MSc degrees in computer science from the University of Hertfordshire. He is currently working toward the PhD degree at the Science and Technology Research Institute of the same university. His research interests include software defect prediction and, in particular, issues with methodology and data quality. His main areas of academic interests include software engineering, machine learning, data quality, and data cleansing.



Steve Counsell received the BSc degree in computing in 1987, the MSc degree in systems analysis in 1988, and the PhD degree from the University of London in 2002. From 1998-2004, he was a lecturer in computer science at Birkbeck, London. He is currently a reader in the School of Information Systems, Computing and Mathematics at Brunel University. His research interests include software metrics, refactoring, software testing, and, more generally, empirical studies of software engineering artifacts.

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