Defect Reduction Planning (using TimeLIME)

Kewen Peng, Tim Menzies, Fellow, IEEE

Abstract—Software comes in releases. An implausible change to software is something that has never been changed in prior releases. When planning how to reduce defects, it is better to use plausible changes, i.e., changes with some precedence in the prior releases.

To demonstrate these points, this paper compares several defect reduction planning tools. LIME is a local sensitivity analysis tool that can report the fewest changes needed to alter the classification of some code module (e.g., from "defective" to "non-defective"). TimeLIME is a new tool, introduced in this paper, that improves LIME by restricting its plans to just those attributes which change the most within a project. In this study, we compared the performance of LIME and TimeLIME and several other defect reduction planning algorithms. The generated plans were assessed via (a) the similarity scores between the proposed code changes and the real code changes made by developers; and (b) the improvement scores seen within projects that followed the plans. For nine project trails, we found that TimeLIME outperformed all other algorithms (in 8 out of 9 trials).

Hence, we strongly recommend using past releases as a source of knowledge for computing fixes for new releases (using TimeLIME). Apart from these specific results, the other lesson from this paper is that our community might be more careful about using off-the-shelf AI tools, without first applying SE knowledge (e.g. that past releases are a good source of knowledge for planning defect reductions). As shown here, once that SE knowledge is applied, this can result in dramatically better reasoning.

Index Terms—Software analytics, Defect Prediction, Defect Reduction, Plausibility Analysis, Interpretable Al

1 INTRODUCTION

"Don't tell me where I am, tell me where to go." – a (very busy) developer

Machine learners generate models. People read models. What learners generate the kind of models that people want to read? If the reader is a busy software developer, then they might not need, or be able to use, complex models. Rather, such a busy developer might instead just want to know the *least* they need to do to achieve the *most* benefits. For example, suppose some AI model has classified a module as "defective". If a developer then asks "what can I do to fix that?" then, ideally, we should, be able to reflect on the model to learn a defect reduction plan; i.e., a small set of actions that reduces the odds of that module being defective.

For many machine learning algorithms, it can be (very) difficult to learn a succinct reduction plan by reflecting on the arcane internal structure of, say, a neural net classifier. To better support busy developers, explanation algorithms like LIME [1] (first presented at KDD'16) can report what attribute changes can alter a classification (e.g., from "defective" to "non-defective"). But classic LIME has a problemit generates surprising and unprecedented plans that had never been seen before in the history of a project. As we show, such unprecedented plans are sub-optimal.

To fix this problem, TimeLIME adds SE knowledge to LIME. We note that software comes in releases and that an implausible change to software is something that has never been changed in prior releases. Hence, we propose the following *TimeLIME tactic*:

When reasoning about changes to a project, it is best to use changes seen in the historical record of that project.

To assess the value of this TimeLIME tactic, we ask:

• **RQ1: Does TimeLIME provide succinct plans?** Classic LIME, proposes changes to dozens of attributes. Time-

LIME, on the other hand, restricts itself to just the most changed attributes. Hence, our plans are easier to apply.

- RQ2: Could developers apply the changes proposed by TimeLIME? Given project information divided into *oldest*, *newer*, and *most recent* data, this paper:
 - Used the *oldest* data to determine what attributes were often changed in a project,
 - Used the *newer* data to build plans using LIME, Time-LIME, and five other planning algorithms;
 - Divided the *most recent* data into:
 - * Those projects that *followed* the plans;
 - * And those that did not. This study found a large overlap (median=80%) between TimeLIME's recommendations and actual actions made by developers.
- **RQ3:** Is TimeLIME better at defect reduction? Time-LIME's plans perform best (compared to classic LIME and four other algorithms).

The rest of this paper is structured as follows. §2 discusses defect prediction, code refactoring, and challenges of using human opinions in SE.§3 introduces some prior works in the field of defect reduction and their methodologies. §4 presents the basic framework of LIME as well as Time-LIME, the new method proposed in this paper. §5 shows our method for ranking different planning methods. §6 describes experiment and the datasets, predictive model, and planners evaluated in this work. §7 and §8 report and discuss the result respectively. The credibility and reliability of our conclusions is discussed by §9. Recent related works are shown in §10, which also declares the major difference that distinguishes the contribution of this paper. Future work and directions are illustrated in §11. Finally, we conclude this work in §12.

1.1 Reproduction Package

All our scripts and data are available on-line¹.

K. Peng and T. Menzies are with the Department of Computer Science, North Carolina State University, Raleigh, USA. E-mail:kpeng@ncsu.edu, timm@ieee.org

^{1.} https://github.com/ai-se/TimeLIME

IEEE TRANSACTIONS ON SOFTWARE ENGINEERING

Metric	Name	Description							
amc	average method complexity	Number of JAVA byte codes							
avg_cc	average McCabe Average	McCabe's cyclomatic complexity seen in class							
ca	afferent couplings	How many other classes use the specific class.							
cam cohosion amongst classes		Summation of number of different types of method parameters in every method divided by a multiplication							
	concision amongst classes	of number of different method parameter types in whole class and number of methods.							
cbm	coupling between methods	otal number of new/redefined methods to which all the inherited methods are coupled							
cbo	coupling between objects	Increased when the methods of one class access services of another.							
ce	efferent couplings	How many other classes is used by the specific class.							
dam	data access	Ratio of private (protected) attributes to total attributes							
dit	depth of inheritance tree	It's defined as the maximum length from the node to the root of the tree							
ic	inheritance coupling	Number of parent classes to which a given class is coupled (includes counts of methods and variables inherited)							
lcom	lack of cohesion in methods	Number of pairs of methods that do not share a reference to an instance variable.							
1		If <i>m</i> , <i>a</i> are the number of methods, attributes in a class number and $\mu(a)$ is the number							
Icom3	another lack of conesion measure	of methods accessing an attribute, then lcom3 = $\left(\left(\frac{1}{a}\sum_{i}^{a}\mu(a_{j})\right) - m\right)/(1-m)$							
loc	lines of code	Total lines of code in this file or package.							
max_cc	Maximum McCabe	Maximum McCabe's cyclomatic complexity seen in class							
mfa	functional abstraction	Number of methods inherited by a class plus number of methods accessible by member methods of the class							
moa	aggregation	Count of the number of data declarations (class fields) whose types are user defined classes							
noc	number of children	Number of direct descendants (subclasses) for each class							
npm	number of public methods	Npm metric simply counts all the methods in a class that are declared as public.							
rfc	response for a class	Number of methods invoked in response to a message to the object.							
wmc	weighted methods per class	A class with more member functions than its peers is considered to be more complex and more error prone.							
defect	defect	Number of bugs which can be transformed into Boolean values for classification.							

TABLE 1: The C-K OO metrics used in defect prediction. The last variable "defect" is the dependent variable.

2 BACKGROUND

2.1 Challenges with Using Human Opinions

This paper is an *algorithmic analysis of historical SE data* where we ran simulations over the historical record of eight software projects. An alternate approach to this algorithmic analysis of historical SE data is to use *qualitative methods*. Qualitative methods rely on surveys of human subject matter experts (e.g., programmers). Much has been learned from such studies of subject matter experts [2]. Nevertheless, in the particular case of large scale defect prediction, we prefer our algorithmic approach, for two reasons:

- *Scalability*: It is hard to scale qualitative investigations of human beliefs to a large number of projects. We mention this since while this paper studies just eight projects, our long-term goal is to develop software analysis methods that applies to hundreds to thousands of projects. While some progress has been seen recently with scaling qualitative methods [3], at the time of this writing, we assert that it is far easier to scale an algorithmic analysis of historical SE data.
- Lack of consensus: multiple studies report that human beliefs in software quality may often be inconsistent and even incorrect. Devanbu et al. have conducted a case study among 564 Microsoft software developers to show that human beliefs on software quality can be quite varied and may not necessarily correspond with actual evidence within current projects [4], [5]. Similar assertions are also made in Passos' paper, where the author reports that conflicting beliefs can be held by different stakeholders of the software development team. There also exist cases that a belief is correct for past projects but not the current work [6]). A more recent study by Shrikanth et al. also reports such much variability of human beliefs about defect prediction [7]. Shrikanth studies 10 beliefs held by software developers about defect prediction, which were initially summarized by Wan et al in 2018 [8]. By measuring the actual support of these beliefs within the project, Shrikanth found that:

- Among over 300,000 changes seen in different opensource projects, only 24% of the projects support all 10 beliefs.
- What is believed the most by developers does not necessarily have the strongest support within projects. For example, a belief that is acknowledged by 35% of the developers has the most support whereas a belief held by 76% of the developers is only ranked 7th out of 10 beliefs.
- As a project grows to mature, the beliefs actually tend to be weakened rather than strengthened.

Not only do practitioners have conflicting beliefs about what causes defects, but we also can see that researchers who have studied many projects also disagree on what factors matter the most to defect reduction. For example, as discussed later in the paper, Alves [9], Shatnawi [10], and Oliveira [11] all offer different models about what matters most for software quality.

In summary, many studies report a significant disconnect between human beliefs and patterns supported by data. Hence, we are nervous about using the opinion of experts' opinions as the "ground truth" to evaluate (e.g.,) defect reduction plans. Accordingly, we use an algorithm analysis since that can use historical SE data to generate the ground truth needed to evaluate a method.

2.2 Defect Prediction

The case study of this paper comes from defect prediction and planning. This section discussed the value of that kind of analysis.

During software development, testing often has some resource limitations. For example, the effort associated with coordinated human effort across a large code base can grow exponentially with the scale of the project [13]. Hence, to effectively manage resources, it is common to match the quality assurance (QA) effort to the perceived criticality and bugginess of the code. Since every decision is associated with a human and resource cost to the developer team, it is impractical and inefficient to distribute equal effort to every component in a software system [14]. Learning defect prediction (using data miners) from static code attributes (like those shown in Table 1) is one very cheap way to "peek" at the code and decide where to spend more QA effort.

Recent results show that software defect predictors are also competitive widely-used automatic methods. Rahman et al. [15] compared (a) static code analysis tools FindBugs, Jlint, and PMD with (b) defect predictors (which they called "statistical defect prediction") built using logistic regression. No significant differences in cost-effectiveness were observed. Given this equivalence, it is significant to note that defect prediction can be quickly adapted to new languages by building lightweight parsers to extract code metrics. The same is not true for static code analyzers - these need extensive modification before they can be used in new languages. Because of this ease of use, and its applicability to many programming languages, defect prediction has been extended many ways including:

- 1) Application of defect prediction methods to locating code with security vulnerabilities [16].
- 2) Understanding the factors that lead to a greater likelihood of defects such as defect prone software components using code metrics (e.g., ratio comment to code, cyclomatic complexity) [17], [18] or process metrics (e.g., recent activity).
- 3) Predicting the location of defects so that appropriate resources may be allocated (e.g., [19])
- 4) Using predictors to proactively fix defects [20]
- 5) Studying defect prediction not only just release-level [21] but also change-level or just-in-time [22].
- 6) Exploring "transfer learning" where predictors from one project are applied to another [23], [24].
- Assessing different learning methods for building predictors [25]. This has led to the development of hyperparameter optimization and better data harvesting tools [26], [27].

This paper extends defect prediction and planning in yet another way: exploring the trade-offs between explanation and planning and the performance of defect prediction models. But beyond the specific scope of this paper, there is nothing in theory stopping the application of this paper to all of the seven areas listed above (and this would be a fruitful area for future research).

2.3 Code refactoring

Code refactoring is an important part of software maintenance. The process is meant to improve the internal quality of software by better structuring the existing code, without changing the external behavior [28], [29]. Such restructuring is assumed to positively affect the software quality by reducing complexity, enhancing maintainability, etc. [30], [31].

Much research has studied the relation between the code refactoring process and internal software static attributes metric and external software quality attributes like maintainability, modifiability, and understandability [32]–[35]. Studies have shown a correlation between external quality attributes and internal quality attributes (such as the OO metrics used in this paper) [36]–[40].

That said, a missing piece of current research is what we call the planning problem. Given that a developer *can* change these metrics in many ways, what should she *actually do*? How do we bridge the gap between "what" to refactor and "why" we need to refactor in the first place? Table 2 shows the effect of various code refactoring methods on our code base. In order to select "what" refactoring method to apply, developers need some mapping from those refactoring to some higher-level goal. The planning algorithms of this paper provide that mapping to the higherlevel goal of "defect reduction in future releases".

3 PRIOR WORK IN PLANNING DEFECT REDUCTION

Over the years, several researchers have proposed various ways to identify appropriate changes on code metrics. This

	AMC	AVG_CC	CA	CAM	CBM	CBO	CE	DAM	DIT	Ŋ	LCOM	LCOM3	LOC	MAX_CC	MFA	MOA	NOC	MAN	RFC	WMC
① inline methods	dec?	inc											dec?	inc?				dec?	dec	dec
2 extract method	inc?	dec											inc?	dec?				inc?	inc	inc
③ extract class	dec	dec		inc?		inc	inc				dec?	dec?	dec	dec?		dec?			dec	dec
④ inline class	inc	inc		dec?		dec	dec				inc?	inc?	inc	inc?		inc?			inc	inc
⑤ move method	dec	inc?											dec	dec?		dec?		dec?		dec
⑥ hide delegate			dec			dec														
⑦ consolidate cond	dec	dec		dec							inc	inc	dec	dec?					inc	inc
(8) polymorphism	dec	dec					inc		inc	inc				dec?			inc			
I flatten cond	dec	dec											dec	dec?						
10 hide method				inc				inc												
(1) simplify para	dec	dec												dec						
(12) factory method	inc	inc											inc	inc?					inc	inc
13 push down method	dec	dec?			dec														dec	dec
(14) encapsulate field				dec							inc	inc	inc							inc
(15) extract subclass	dec	dec	inc	inc?			inc				dec?	dec?	dec	dec?	dec		inc		dec	dec
(16) inline subclass	inc	inc	dec	dec?			dec				inc?	inc?	inc	inc?	inc		dec		inc	inc

TABLE 2: Code refactoring methods. Taken from [12]. In this table "dec" and "inc" are short hand for decrease and increase (respectively). These cell values were determined as follows. For each of the methods in column one, ten times, we applied that refactoring method to some randomly selected portion of the code base used in this study. The measurements listed in the columns headers were collected before and after that change. Empty cells indicate no influence on the feature. Cells marked with "?" means we observed that some examples have the change while others do not.

section will illustrate 4 methods that rely on either *outlier statistics* or *cluster deltas*.

Outlier statistics: The general principle underlying outlier statistics methods is that in the distribution of values for each code metric, there are some extremely large/small values that are associated with greater defect proneness. Therefore, by changing those metrics to not have such outlier values, the code base may be found fewer bugs. This paper presents 3 outlier statistics methods and the major distinction among them is their different ways to identify the threshold for outlier values. In the following text, the methods of Alves et al., Oliveira et al, Shatnawi are based on outlier statistics.

Cluster deltas is a framework for learning conjunctions of rules that need to be applied to the code metrics simultaneously. Unlike outlier statistics, which merely studies the statistical distribution of code metrics, cluster deltas is a supervised learner that take account of whether the code base is defective. In the following text, Krishna's XTREE method uses cluster deltas to learn association rules concerning about when and where to apply a code change.

These two approaches are discussed below. Before doing that, we first digress to make the point that *none* of the following methods can be effective unless:

- 1) It can be shown that programmers can apply the suggestions made in these plans;
- 2) It can also be shown that when programmers apply the suggestions, they do not inadvertently add other changes that reduce (or remove) the effectiveness of these plans.

Later in this paper, we show evidence that these points 1,2 are actually achievable – see §8.

3.1 Alves. 2010

Alves et al. [9] offers an unsupervised approach that learns from the statistical distribution and scale of OO metrics. At the beginning, Alves' method will weight each metric value according to the lines of code (LOC in Table 1) of its code class. The weighted metric values will then be normalized by the total sum of weights and sorted in an ascending order. Note that the sorted result is just equivalent to a cumulative probability function where x-axis stands for the weight percentage from 0 to 100% and y-axis the metric scale.

After that, a threshold percentage will be customized (Alves et al. recommends 70%) to identify normal metric values against abnormal metric values. For example, a threshold of 70% will identify the value for each metric where 70% of the classes fall below. The intuition behind this is straightforward: they believe that a code class with outlier metric values that exceed 70% of its peers is more likely to be found bugs.

When we implemented the Alves' method in our experiment, we augmented the original implementation by also studying the correlation between the code metrics and the defect state of the class. By fitting each dependent variable and the independent variable with a univariate logistic regression classifier:

• we were able to reject metrics that are poor indicators of defects (here we define "poor" as a logistic regression with *p*-value > 0.05).

- For those metrics that survived from the rejection, the planner will identify the normal range according to the threshold, i.e., [0, 70%] for each metric.
- Finally, during the planning process, any "survived" metric exceeding the threshold value will be proposed to reduce its value to the normal range.

3.2 Shatnawi, 2010

Shatnawi [10] in 2010 provided an alternative to Alves' method by using VARL (Value of Acceptable Risk Level) to compute the outlier threshold. Initially proposed by Bender [41] in 1999 in his epidemiological studies, the VARL function is a supervised learner that uses the interpretation of the univariate logistic regression model to derive the threshold for an acceptable risk level given by a probability p_0 (i.e., $p_0 = 0.05$). That is to say, the VARL believes that the probability p_0 of an event is less than 0.05 of the value of the dependent variable is smaller than VARL. The VARL function is as follow:

$$VARL = \frac{1}{\beta} (log(\frac{p_0}{1-p_0}) - \alpha)$$

Here, α is the intercept of the logistic regression, β is the coefficient of the logistic regression, and p_0 is the acceptable risk probability as described above.

Similar to our procedure of implementing Alves' method, we ruled out metrics with *p*-value > 0.05, and computed the VARL for the remaining metrics. We define the proposed plan for each metric as [0, VARL], which means a metric value exceeding VARL will be recommended a reduction by the planner.

3.3 Oliveira, 2014

Oliveira et al. [11] approach an totally different threshold definition than the previous 2 methods. Instead of deriving an absolute threshold like Alves et al. and Shatnawi did, Oliveira et al. choose to use the *relative threshold*, which indicates the percentage of classes the the upper bound (threshold) shall be applied to. The general format of their defect reduction rules is as follow:

p% of the classes must have $M \leq K$

Here, M is the code metric; K is the threshold value for the corresponding metric; p% is the minimum percentage of code classed that are required to follow the restriction specified above.

In order to compute the pair of values (p, K) for each metric M, Oliveira defines 3 functions: **Compliance**(p, k), **Penalty1**(p, k), and **Panelty2**(p, k). The **Compliance** method reports the percentage of classes that follow the rule defined by each pair of values (p, K). The **Penalty1** penalizes the model if the compliance rate is lower than a constant percentage (i.e., 90%). **Penalty2** computes the distance between k and the median of the preset *Tail*-th percentile for each metric (Oliveira et al. suggest 90-th percentile). Summing up the 2 penalty values to obtain the total penalty, the method chooses the pair of values (p, K) with the lowest total penalty where a tie will be broken by choosing the highest p and the lowest k.



Fig. 1: XTREE plan generation. From Krishna et al. [42]. An example has fallen down to the current branch where the probability of defects is 1.00. The nearby desired branch predicts a 0.00 probability of defects. XTREE's plans are the delta Δ between the branches.

3.4 XTREE, 2020

Earlier in 2020, Krishna [42] proposed XTREE, a novel defect reduction planning method that does not rely on outlier statistics. The XTREE planner uses frequent pattern mining, decision trees, and a walk traversal algorithm.

With the pattern mining, XTREE attempts to find what code metrics usually change together by applying an association rule learner on historical data. Since metrics in Table 1 are continuous, XTREE will first discretize the values into intervals using Fayyad-Irani. Then a FP-growth algorithm [43] is used to mine frequent itemsets (in our experimentation XTREE uses $minSupport = 5\% \times total_size$).

The returned maximal frequent itemsets are used to construct a decision tree. After that, in the third part, the plans will be generated by traversing the decision tree to seek for the closest branch with highest improvement in the probability of the non-defective label. An example of the traversal procedure is illustrated in the Figure 1. Once the **current branch** is found, the plan will be the Δ from the current branch to a nearby **desired branch** with lower probability of defects.



Fig. 2: An example of output generated by Table 3 when applied to the data sets of the form of Table 1. The y-axis shows the feature name and the confidence interval during which the explanation stays effective. The x-axis indicates the importance weight of each attribute. The prediction label of this instance is 1 (defective), and the weights show how each feature contributes to the prediction. A positive weight means the feature encourages the classifier to predict the instance as a positive label (defective), and vice versa for the negative weight. Larger weights indicate greater feature importance in terms of the prediction value based on that feature weighted by a similarity kernel.

4 New Methods for Planning Defect Reduction (LIME and TimeLIME)

4.1 LIME

One of the starting points of this research was the realization that the LIME algorithm, first published at KDD'16 [1] could be applied to defect reduction planning. The internal framework of LIME is depicted in Table 3. In summary, given an instance I of class X, LIME conducts a sensitivity analysis in the neighborhood around I to determine what could change the class from X to Y. Using the synthetic data generated around I, LIME can get the classification/regression result from any black-box learner, which will then be used to fit a linear model describes the local region. The parameters of the fitted linear model are then reported as a way to understand how changes in values can adjust the classification; e.g., see Figure 2.

This paper utilizes LIME and its capability in interpretation to generate defect reduction plans. If a black-box model can predict defects accurately, then it might be "knowledgeable" enough to provide more informative plans than a subject matter experts can provide. The key question is,



TABLE 3: Inside LIME. From [1]. The feature importance weights are passed to Algorithm 1 and 2, as later elaborated in §6.3. For a sample of the output feature importance weights, see Figure 2.

therefore, how could we access the knowledge owned by a black-box model. In this paper, we imported LIME as the core component of our defect reduction algorithm as we also leverage other software domain knowledge to help LIME restrict the proposed plans in an effective fashion.

Sometimes, we are asked why we are basing our approach on LIME and not other other tools that explain how to change attributes in order the change the classification of an instance. To say the least, there are very many alternate algorithms. A recent survey by Mueller et al. summarized various kinds of change-explanation generation tools. [44]. Mentioned in their study, Mueller et.al report that this literature is truly vast. Consequently, there are many alternatives to LIME including the abductive framework of Menzies et al. [45] or ANCHORS [46] (which is another changeexplanation algorithm generated by the same team that created LIME).

We based our work on LIME, for several reasons. Firstly, LIME scales to large problems. Much recent work has results in methods to scale data mining to very large data sets. Since LIME is based on data mining, then LIME can use those scalability results in order to generate explanations for very large problems.

Secondly, and this is more of a low-level systems reason, alternatives to LIME such as ANCHORS assume categorical or discrete features. Our data has continuous classes which could be binarized into two discrete classes– but only at the cost of losing the information about local gradients. Hence, at least for now, we explore LIME (and will explore ANCHORS in future work).

Lastly, LIME is a widely-cited algorithm. At the time of this writing, LIME has received over 3,000 citations since it was published in 2016. Hence, methods used to improve LIME could also be useful for a wide range of other research tasks. This paper proposes precedence plausibility as a way to improve LIME.



Fig. 3: TimeLIME: overview of the algorithm, plus the *K*-test evaluation rig. Note that for evaluating other benchmark algorithms, the area bounded by the dotted line will be replaced by the corresponding algorithm. For further details on TimeLIME, see Algorithm 2.

4.2 TimeLIME

TimeLIME extends LIME by restricting the generated recommendations to the attributes which were seen to be frequently modified within the history of a software projects. Figure 3 offers a graphical overview of this system.

TimeLIME evolved out of comments we heard at workshop on "Actionable Analytics" at ASE'15 [47]. There, business users complained about analytic models saying that rather applying a black-box data mining algorithm, they preferred an approach with a seemingly intuitive appeal. Since software engineers are the target audience of analytics in SE, it is crucial to ensure the proposed recommendations are valued by them. Chen et al. say the term "actionable" can be defined as a combination of "comprehensible" and "operational" [21]. But how to assess "operational"?

In this paper we make the following assumption about "operational": a proposed change to the code is plausible if it has occurred before. That is, in this work, we claim a plan is the most operational when it has the most precedence in the history log of the project.

Using this assumption, we can generate operational analytics by:

- Looking at two releases of a project and report the attributes that have changed between them;
- Next, when generating plans, we only used those attributes that have the most changes.

After conducting a survey on 92 controlled experiments published in 12 major software engineering journals, Kampenes et al. [48] argues that in SE, size change can be measured via Hedge's g value [49]:

$$g = (M_1 - M_2)/(S_{pooled})$$
 (1)

Here, M_1 and M_2 are the means of an attribute in two consecutive releases and S_{pooled} comes from 2. This expression is the pooled and weighted standard deviation (n and s denote the sample size and the standard deviation respectively).

$$S_{pooled} = \sqrt{((n_1 - 1)s_1^2 + (n_2 - 1)s_2^2)/(n_1 + n_2 - 2)}$$
 (2)

For the details on how Equation 2 was applied, see §6.3.

Furthermore, in order to ensure the precedence of plans generated by TimeLIME, we will only allow changes that have actually been "seen" in the record before. When a plan is generated given a certain set of actionable features (we denote such set as the candidate pool), we look it up in the historical records of the project. If there exist such records where the exact changes proposed by TimeLIME had once been made by developers, we will return the plan. Otherwise, the algorithm will go back to the previous step and start generating a different plan using the same candidate pool. The plan generating process is a greedy attempt: Only if all the combinations of M candidates fail, we will try to make a plan using (M - 1) candidates. In theory, TimeLIME might generate conflicting plans (but in our logs, this occurs very rarely). Nevertheless, to handle that situation, we recommend that if two plans are conflicting, users should adopt the plan with greatest support (i.e. most frequent in the historical log).

The reason why TimeLIME uses association rule mining to ensure the precedence of plans is that TimeLIME is designed to provide not **optimal** solutions, but **achievable** and **maintainable** solutions. Different development teams have different capacity of reducing defects. If a planner is learning from the historical activities of developers (which contains examples of good defect reductions) to generate plans, then better plans could be generated by learning from better development teams. Thus, it is impossible to generate an optimal plan since no one can assert whether the current plan is optimal. Alternatively, it is viable to generate a maintainable plan: if the current team has been working on reducing defects in a satisfactory (may not be best) level, TimeLIME would help to maintain the defect reduction quality that has been achieved so far.

5 Assessing Planners: The K-TEST

This paper claims that plans from TimeLIME planner (that focus on attributes with a history of most change) outperform those generated from LIME, XTREE, Alves, Shatnaw, and Oliveira. To defend that claim, we need some way to assess different planning methods.

There's an expression in Latin, *post hoc ergo propter hoc*, which means "after this, then because of this". This expression refers to the logical fallacy that "if event B follows event A, then event A must be the cause of event B". The assertion is obviously flawed since other events could be the true trigger of event B. This is why, in this study, we need to carefully evaluate the effectiveness of plans to see if knowledge learned from past code change records actually helps make plans on future code changes.

To address this concern, we use Krishna's *K*-test [42]. The *K*-test uses historical data from multiple software releases to compare the effectiveness of different plans P_1, P_2, \ldots . The test is a kind of simulation study that assumes developers were told about a plan at some prior time. Given project information divided into *oldest, newer,* and *most recent,* we will use the *oldest* data to determine what attributes where often changed in a project. Then, using the *newer* data, we will build plans using LIME, TimeLIME, XTREE, Alves, Shatnaw, and Oliveira. Finally, we will divide the changes between the *newer* and the *most recent* into the changes that *overlap* with the plans, and those that do not.

More precisely, we use consecutive releases x, y, z of some software system. These releases are required to contain named regions of code C_1, C_2 , etc. that can be found in releases x, y, z. For example, C_i could be an object-oriented class or a function or a file that is found in all releases. The *K*-test then assumes that there exists a quality measure Q that reports the value of the regions of named code in different releases. In this study, we will use NDPV (*Number* of *Defects in Previous Version*) as the quality measure, which is described later in §6.4. Some method is then applied that uses Q to reflect on the releases x, y in order to infer a plan P_i for improving release z^2 .

Given the above, the *K*-test collects following quantities to address our claims made in introduction:

• **RQ1: Smaller**: To measure the succinctness of plans, we collect the number of changes within the plan proposed for each code *C_i* in release *y*.

- RQ2: Ready to apply: To measure how likely a plan can be realized by developers, we compute J_{y,z} = Δ_{y,z} ∩ P_i: the overlap between the proposed plan and the code changes.
- **RQ3:** Better: To measure which planner is better at reducing defects, we collect $Q_z Q_y$: i.e., the change in the number of bugs of the named code C_i between releases y, z. Then, we weight the change $Q_z Q_y$ by $J_{y,z}$. The intuition is that the planner cannot get credit in a bug-reducing code file if its plan shares little or none similarity with the actual actions done by developers.

The *K*-test defines *better* plans as follows:

DEFINITION: Plan P_i is "better" that plan P_j if, in release z, P_i is associated with most quality improvements.

That is, increasing the size of the overlap of the proposed plan is associated with increasing quality in release *z*; i.e.,

$$(Q_z - Q_y) \propto |J_{y,z}|$$

For our purposes, the *K*-test procedure in this paper consists of three steps:

- Train a defect reduction planner on version *x*.
- Use trained planner to generate plans with the aim of fixing bugs reported in version *y*. In this step, classical LIME planner and TimeLIME planner will utilize the explanations from the explainer and TimeLIME, in addition, will use the historical data analysis to generate plans.
- On the same set of files that are reported buggy in version y, we measure $J_{y,z}$, the overlap score of each plan and the changes in the version z, using the Jaccard similarity function. We also record $Q_z Q_y$, the change in the number of bugs between the version y and version z.

For each instance, we compare the extent of overlap between the recommended plan P_i generated by the planner and the actual developer action in the next release as $\Delta_{y,z}$ using the Jaccard similarity coefficient.

$$J_{y,z}(P_i, \Delta_{y,z}) = (P_i \cap \Delta_{y,z})/(P_i \cup \Delta_{y,z})$$
(3)

Then we convert the coefficient into percentage as our overlap score. As an example shown in Table 4, the overlap score is

 $2/4 \times 100\% = 50\%$

	AMC	LOC	LCOM	CBO
Current release y	0.2	0.1	0.9	0.5
P_i for release z	no change	[0, 0.1]	[0, 0.9]	no change
Next release z	0.2	0.3	0.3	0.2
Match?	у	n	у	n
Map to Table 10	TN	FP	TP	FN

TABLE 4: A contrived example: computing similarity score using the Jaccard function from Equation (3). Plans that match the developer actions are marked gray.

Formally speaking, the *K*-test is *not* a deterministic statement that some plan will necessarily improve quality in some future release of a project. Such deterministic causality is a precisely defined concept with the property that a single counterexample can refute the causal claim [51]. The *K*-test does <u>not</u> make such statements.

Instead, the *K*-test is a statement of historical observation. Plans that are "better" (as defined above) are those

^{2.} Note the connection here to temporal validation in machine learning [50]. In the K-test, no knowledge of the final release z is used to generate the plans.

which, in the historical log, have been associated with increased values on some quality measure. Hence, they have some likelihood (but no certainty) that they will do so for future projects.

6 EXPERIMENTAL METHODS

The experiment reports the performance of TimeLIME and other state-of-the-art works by comparing the quality of plans recommended by each method.

Firstly, we use an over-sampling tool called SMOTE [52] to transform the imbalanced datasets in which defective instances may only take a small ratio of the population. This was needed since, in many of the prior papers that explored our data, researchers warn that small target classes made it harder to build predictors [53].

Secondly, as discussed above, we train the predictor P and explainer E on data of version x. Then in version y we use the explainer to generate explanations *only* on those data that are reported as buggy. We also use the predictor P to determine whether we should provide recommendation plans to the instance.

Then we measure the overlap score of our recommended plan and the actual change on the same file in version z. To do this, only select instances that are defective and whose file name has appeared in all releases of data to be instances in need of plans.

The above steps are applied for each benchmark method as well as the TimeLIME planner proposed by this paper. The visualization of the experimental rig is shown in Figure 3. In the classical LIME planner, we use the simple strategy which is to change as many features as it can in order to reduce the number of bugs. On the other hand, for Time-LIME, we first input historical data from the older release to compute the variance of each feature. Then we selected the top-M features with the largest variance as *precedented* features, meaning any recommendation on other features will be rebutted. After getting recommended plans from both planners, we assess the performance of two planners using the overlap score as described in §6.4.

Note that the parameter M can be user-specified and the features may vary with respect to different projects and the releases used as historical data. Here we set the default value of M to be 5, which means only 25% of all twenty features can be mutated. Our results from experiments suggest that M = 5 is a useful default setting. Future work shall explore and compare other values of M.

6.1 Data

To empirically evaluate classical LIME vs TimeLIME, we use the standard datasets and measures widely used in defect prediction. In this paper, we selected 8 datasets from the publicly available SEACRAFT project [54] collected by Jureczko et al. for open-source JAVA systems (http://tiny. cc/defects). These datasets keep the logs of past defects as shown in Table 5 and summarize software components using the CK code metrics as shown in Table 1. Note that all the metrics are numerical and can be automatically collected for different systems [55]. The definition and nature of each attribute in the metrics is elaborated by prior researchers Jureczko and Madeyski [56], [57]. Another reason this paper selects these 8 datasets is that they all contain at least 3 consecutive releases, which is required by the evaluation measure described in §5. Since Camel dataset contains 4 consecutive releases, the experiment has 9 trials in total.

6.2 Learner

While other benchmark algorithms don't need the predictive learn within their model, LIME does require the user to pass in the customized learner, which can be used to generate explanations. Since the goal of this paper is to examine the performance of the defect reduction tools rather than the predictive model, this paper takes one classifier to apply the explanation algorithm on.

Our choice of classifier is guided by the Ghotra et al. [58] study that explored 30 classification techniques for defect prediction. They found that all the classifiers they explored fell into four groups and that Random Forest classifiers were to be found in their top-ranked group.

A Random Forest classifier is an ensemble learner that fits a number of decision tree classifiers on different subsamples of the dataset and generates predictions via average voting from all the classifiers [59]. It is impossible to visualize a fitted Random Forest classifier as a finite set of rules and conditions due to the voting process. Therefore, Random Forest classifier is considered a non-interpretable model. Hence, it is a suitable choice for this study.

6.3 Planners

This section discusses the internals of our planners, including a RandomWalk planner (which we use to compare our results against a baseline random guesser).

Using LIME, we generate plans to reduce classifications. We use the default parameter setting of LIME, which is 5000 samples around the instance neighborhood, and the

	Training	Testing	Validation	No. of	No. of	No. of bugs	No. of bugs	No. of
Dataset	(oldest)	(newer)	(most recent)	files	matched files	in testing set	in validation set	bugs reduced
Jedit	4.0	4.1	4.2	367	78	216	74	142
Camel1	1.0	1.2	1.4	872	210	508	247	261
Camel2	1.2	1.4	1.6	965	144	334	316	18
log4j	1.0	1.1	1.2	205	35	83	120	-37
Xalan	2.5	2.6	2.7	885	385	529	381	148
Ant	1.5	1.6	1.7	745	91	183	163	20
Velocity	1.4	1.5	1.6	229	138	321	144	177
Poi	1.5	2.5	3.0	442	247	495	366	129
Synapse	1.0	1.1	1.2	256	58	97	65	32

TABLE 5: Defect datasets used in this paper. Each columns represent a different dataset. The last row shows the total number of bugs reduced among the same files between the testing release and the validation release. Note that a negative value in this column indicates that the validation release contains more bugs than the previous one.



Data: explanation <i>e</i> // the explanation from Table 3							
Result: A tuple consisting of intervals of values v'							
begin							
$w, v \leftarrow e / /$ split weights w and value intervals v							
from e							
$i \leftarrow 0$							
while $i \leq sizeof(w)$ do							
if $w[i] \ge 0$ then							
$v'[i] \leftarrow flip(v[i])$							
else							
$v'[i] \leftarrow v[i] / / \text{ do not propose a change on}$							
this feature							
$\lfloor i \leftarrow i + 1$							
return v'							

Algorithm 2: TimeLIME Planner

Data: explanation *e* from Table 3, precedence parameter M, previous release x, current release y**Result:** A tuple consisting of intervals of values v'begin $w, v \leftarrow e$ // split weights w and value intervals vfrom e $M \leftarrow 5$ // the default parameter M is 5 meaning at most 5 features can be changed in the resulting plan $g \leftarrow \text{hedge}(x, y) / / \text{defined in §4.2}$ precedented $\leftarrow sorted(g)[0:M]$ $i \leftarrow 0$ $pool \leftarrow v$ while $i \leq sizeof(w)$ do if $w[i] \ge 0$ and $i \in precedented$ then $pool[i] \leftarrow flip(v[i])$ else continue // do not propose a change $i \leftarrow i + 1$ $v' \leftarrow findSupport(pool, precedented, x, y)$ // as described in Algorithm 3 return v

entropy-based discretizer. The explanation object return by a LIME explainer is a tuple in which each element contains the feature name and the corresponding feature importance. It also provides a discretized interval indicating the range of values during which the feature will maintain the same effect to the prediction result. As described in Algorithm 1, the simple planner based on the classical LIME will recommend changes on all features that contribute to the defective prediction. Algorithm 2 shows the TimeLIME planner, which utilizes Algorithm 3 to ensure that the proposed plan must be precedented in he historical records. Each planner uses feature ranges generated by flipping the discretized interval relative to the feature value range [0, 1].

Also, just to compare, we use a planner named RandomWalk as a "straw-man" baseline algorithm. This planner, as shown in Algorithm 4, assigns random recommendations to each variable stochastically.

One final note: to make the comparisons fair, in our experiment setting, we set the number of changed features as the same as the TimeLIME planner for comparison purpose.

Algorithm 3: findSupport
Data: candidates changes on each single feature pool,
precedented features <i>precedented</i> , previous
release x , current release y
Result: The proposed plan for the instance
begin
$max \leftarrow 0 / /$ initialize the max support
$itemsets \leftarrow (x, y) / / \text{get records of actual changes}$
$M \leftarrow 5$ // number of changes allowed
while $M > 0$ do
$plans \leftarrow generatePlans(pool, M, precedented)$
$plan, max \leftarrow findMax(plans, itemsets) //$
return the plan with max support in the
historical records
if $max > 0$ then
$\mid break$
else
$ \qquad \qquad$
return <i>plan</i>

Algorithm 4: RandomWalk

Data: standardized code instance to be explained								
c_i number of features to be mutated \hat{n}								
Result: A tuple consisting of intervals of values <i>v</i>								
begin								
$pool \leftarrow random.sample(20, n) // randomly choose$								
n out to 20 features								
$i \leftarrow 0$								
while $i \leq sizeof(c)$ do								
if pin pool then								
$(a, b) \leftarrow sorted(rand(1), rand(1)) / generate$								
a random interval within the range [0, 1].								
$v[i] \leftarrow (a, b) / / apply the random interval.$								
else								
$_$ return v								

6.4 Performance Criteria

The two performance criteria in this experiment, as described in the §5, are the overlap score of individual plans and the number of bugs reduced/added in the next release of the project. The function used for computing the overlap score is the Jaccard similarity function in Eq. 3, and the other criterion is measured by the metric NDPV (*Number of Defects in Previous Version*), which returns the number of bugs fixed (or added) in a given file during the development of the previous release. The nature of NDPV and similar metrics have been evaluated by plentiful researchers [60]–[63].

To further evaluate the second criterion, we chose to use a weighted sum function to compute the net gain of each planner. The weighted sum function in Eq. (4) weights the NDPV by the overlap score of the plan.

$$S = \sum s_i * n_i \tag{4}$$

In the experiment, each plan p_i from the all N plans returns an overlap score s_i and a NDPV number n_i (positive number indicates bugs reduced, negative number indicates bugs added). Then we multiply the NDPV n_i by s_i to compute the weighted improvement score S. Note that the *larger* s_i indicate the *greater* overlap. Fig. 4 shows the tendency of the



• Plan D: High similarity to a bug-adding action

Fig. 4: Visualized tendency of the Eq. 4. The x-axis shows the NDPV n_i and the y-axis shows the similarity score s_i .

weighted improvement score S with respect to s_i and n_i . Generally speaking, we will reward plans who are similar to a bug-reducing change and penalize those plans who are similar to a bug-introducing change. In the case where the plan is very dissimilar to a change (whether it is bug-reducing or not), we assign a trivial score to the plan since it shares little overlap with the actual change, which makes it impossible for us to simulate the potential consequence of applying such plan.

Additionally, given that the total number of bugs varies from each project as shown in Table 5, a project with more bugs reduced in the validation dataset will expect the planner to score more than the planner whose validation dataset has fewer bugs reduced so that their performance can be considered proportionally similar. For example, project A has NDPV = 100 in release y and another project B has NDPV = 10 in its next release y. Assume one would like to observe similar performance of a planner on these 2 projects, it won't make any sense if the planner gains the same score in both projects. From this perspective, we scale the final score S in Eq. 4 by the sum of NDPV within the project to get the scaled score S_{scaled} .

$$S_{scaled} = \frac{\sum_{i}^{N} s_i * n_i}{\sum_{i}^{N} n_i} \tag{5}$$

7 RESULTS

7.1 RQ1: Does TimeLIME provide succinct plans?

Figure 5 reports the mean size of plans across all instances in release z. We note that:

- RandomWalk method's plans are so large since this planner does not use information from the domain to constraint its results.
- TimeLIME generates much smaller plans compared to many other planners including classical LIME.
- The only planner that consistently produces smaller plans in the Shatnawi method but, as seen in the <u>RQ3</u> results (below), the Shatnawi obtains performance that is far worse than TimeLIME.



Fig. 5: RQ1 results: The mean size of TimeLIME's plans (across all instances in release z) is often much smaller than LIME. Y-axis shows the number of features changed by recommended plans.

Note that since TimeLIME in the experiment restricts plans to the top 5 features with highest Hedge's *g* scores, the size of an TimeLIME plan will never be more than 5. However, as shown in the figure, the average size of TimeLIME plans is always smaller than 5. This implies that the original code refactoring plans, proposed by the classical LIME planner, do contain unprecedented changes which then get rejected by the TimeLIME planner. In summary, for **RQ1**, we say:

Answer 1: Several planners, including LIME, generate plans that are far larger than those found by TimeLIME. And the only planner that always generates smaller plans has much worse performance.

7.2 RQ2: Could developers apply the changes proposed by TimeLIME?

We answer this question in two ways. Firstly, we assess "can developers map our plans onto known refactoring actions?" (and for our definition of "refactoring actions", please see Table 2). Table 6 shows those mappings. While things get somewhat complicated in two cases (Log4j and Velocity) it is encouraging to note that in $\frac{9}{11}$ cases, the number of refactoring actions is *less than* the number of changes recommended in TimeLIME's plans. Concrete examples of how to apply these mapped plans can be found on our homepage. ³.

For a second way to answer this question, we use our historical data. We posit the scenario that developers were told of our plans in the *current* release, and then we check the *later* release to see if the plans we proposed were actually recommended. As shown in the rest of this section, we say our plans are *feasible* since there is evidence indicating developers could actually apply those changes.

^{3.} https://github.com/ai-se/TimeLIME/blob/master/README.md#examples

IEEE TRANSACTIONS ON SOFTWARE ENGINEERING

	Jedit	Camel1	Ca	mel2	Log4j	Xa	lan	Ant	Velocity	Poi	Synapse
TimeLIME Plans	- lcom3	- lcom3	+ lcom3	- moa	- max_cc	- cbo	- loc	- rfc	- cbm	+ rfc	+ avg_cc
	- moa	+ dam	+ moa	- ce	- avg_cc	+ loc	+ cam	- ce	- mfa	+ amc	- cbm
	- avg_cc	+ cam	- ce	- lcom3	-wmc	+ amc	+ cbo	- npm	- amc	- dam	- mfa
	- max_cc	- ic	+ rfc	- rfc	- npm	- cam	+ ce	-wmc		+ loc	- cam
	+ cam										
Refactoring Methods	Q, 3, 15	3, 15	4,16	3, 5, 15	3, 5, 8, 9 (1), (13, (15)	(), Ø	3	1), (5), (13)	(1), (7), (8), (9) (11), (13), (15)	@, 12, 16	4, 13

TABLE 6: Using Table 2, developers can map TimeLIME's plans onto some simple refactoring methods to achieve the desired changes in code metrics. The "+" and "-" indicate increase and decrease respectively. Note that while various plans are provided within each project, in this table we only show the most frequent plan(s).

Table 7 and Table 8 comments on how often developers are willing to perform the plans suggested by different planners. Both tables are generated using the *K*-test procedure described above. Each cell in Table 7 shows the median value of the $J_{y,z}$ overlap score measured from Eq. 3 in §5 for all instances within the projects. In addition, in order to explore the robustness of our approach, we added 2 variants of TimeLIME planners, each embedded with a different predictor. While the original TimeLIME planner uses a random forest classifier internally, the 2 variants use multi-layer perceptron (MLP) and support vector machine (SVM) respectively.

Table 8 shows the interquartile range (IQR) of all overlap scores quantile among all plans generated by planners. With similar median scores, a smaller IQR means the planner is more stable and robust. It is noteworthy that the Random Planner always obtains very small IQRs in all project. This is because plans generated by Algorithm 4 are equivalently bad as indicated from the median scores. On contrary, TimeLIME has similarly small IQRs while maintaining the highest median scores in all project, which means it prevails other planners in terms of providing plans that better resemble developers' choices. In summary:

- Unsurprisingly, Random Planner has the lowest similarity scores in all projects.
- The 4 prior works (XTREE, Alves, Shatnawi, and Oliveira) are equivalently good.
- Different projects have very different baselines for the similarity evaluation. For example, within Xalan, every planner except Random obtains a relatively high scores whereas they perform equally poorly in the Ant project.
- All TimeLIME planners have obtained the highest score in every project with a relatively low IQR scores. This means the performance of TimeLIME, regardless of the type of the embedded classifier, is good and robust in terms of similarity to actual actions.
- The classical LIME planner has a volatile performance: It is either performs best or worst. In other words, compared to TimeLIME, we cannot recommend that procedure for practical purposes.

In summary, we answer **RQ2** as follows:

Answer 2: We find a large overlap between Time-LIME's recommendations and the possible actions (Table 6) and observed actions (Table 7, Table 8) of developers.

Hence we say developers would be able to apply Time-LIME's recommendations.

	Random	Oliveira	Shatnawi	Alves	XTREE	LIME	TimeLIME	TimeLIME_MLP	TimeLIME_SVM
Jedit	30	35	35	30	40	35	90	85	90
Camel1	55	63	65	60	67.5	85	95	95	95
Camel2	45	55	55	50	50	40	95	95	95
Log4j	45	50	45	40	50	35	75	80	75
Xalan	70	85	90	80	75	85	100	100	100
Ant	30	35	35	35	35	70	85	85	85
Velocity	60	75	75	55	65	50	100	95	100
Poi	35	40	35	40	45	45	75	75	75
Synapse	35	43	45	40	50	73	75	75	75

 TABLE 7: RQ2 results: Median overlap scores in percentage: larger scores are better, marked in darker color.

	Random	Oliveira	Shatnawi	Alves	XTREE	LIME	TimeLIME	TimeLIME_MLP	TimeLIME_SVM
Jedit	30	20	20	15	15	50	15	19	15
Camel1	45	50	55	45	40	20	15	20	15
Camel2	40	50	46	45	50	26	10	25	10
Log4j	28	40	50	38	35	22	25	10	25
Xalan	30	45	35	30	20	25	10	23	5
Ant	18	20	20	15	23	50	22	22	20
Velocity	34	50	45	35	39	20	10	10	10
Poi	18	20	20	20	20	15	20	20	20
Synapse	25	35	44	39	34	39	30	35	34

TABLE 8: RQ2 results: IQR overlap scores: for the same median scores, smaller IQRs are better, marked in darker color.

7.3 RQ3: Is TimeLIME better at defect reduction?

As discussed earlier, better plans in defect reduction field are believed to be those that are (a) easier to apply while (b) maintaining the effectiveness in reducing bugs. The first criterion has already been met. As seen there, the plans made by TimeLIME are much smaller, hence easier to apply, than the other methods studied here. Also, as seen above, the plans from TimeLIME correspond well to the known actions of developers.

The visualized result in Table 9 shows that 3 variants of TimeLIME planners have obtained highest average S_{scaled} scores in most of the projects (8 out of 9).

The overall result is very clear:

Answer 3: The changes proposed by TimeLIME are associated with a much larger reduction in defects than classic LIME and other benchmark algorithms.

	Random	Oliveira	Shatnawi	Alves	XTREE	LIME	TimeLIME	TimeLIME_MLP	TimeLIME_SVM
Jedit	33	34	19	31	40	43	86	84	86
Camel1	45	55	16	50	62	77	84	85	86
Camel2	35	42	13	38	35	63	65	74	65
Log4j	33	44	23	39	42	43	78	77	75
Xalan	60	81	60	74	67	85	96	95	95
Ant	49	58	18	47	62	99	94	97	92
Velocity	50	58	18	47	52	48	91	91	91
Poi	38	50	0	50	54	43	75	77	74
Synapse	37	41	1	36	41	60	71	65	66

TABLE 9: RQ3 results: Improvement percentage perproject: the higher the better. Best and worst planner ineach project are marked in dark and light respectively.

	Actual: change	Actual: no / different change
TimeLIME: change	TP	FP
TimeLIME: don't change	FN	TN

TABLE 10: Each change/no-change proposed in a plan will be categorized into one of the 4 kinds according to the actual value in the *most recent release*. Example can be found in Table. 4

8 DISCUSSION

A potential major objection to all the above could be that the planning process, as we described so far, may be inefficient due to:

- Developers may not be able to implement our plans;
- Even if developers could implement the plans, they might inadvertently make other changes that negate the improvements suggested in the plans.

In this section, we will discuss the practicality of our approach concerning these 2 issues, followed by another section of other, less pressing, threats to validity. To that end, we extended our measurement of similarity between the proposed plans and the actual actions. We further label each change in a single plan into one of the following 4 categories:

- True Positive: suggests same change as seen later;
- True Negative: suggests no change, and no change later;
- False Positive: suggests a change which is not seen later;
- False Negative: suggests no change, but some other change is found later.

We also calculate the precision and recall as defined in Table 10:

- **Precision = TP/(TP+FP)**: Among all changes proposed by a planner, how many of them are found in the next release?
- **Recall = TP/(TP+FN)**: For changes found in the next release, how many are the same as the planner's?

Unanticipated changes that are not recommended by a planner are marked as **FN**. Plans proposed by TimeLIME but not happened in the next release are marked as **FP**. A higher **precision** means more of TimeLIME's plans are undertaken, and a higher **recall** means there are fewer unanticipated changes in the next release. If it is TimeLIME rather than unanticipated changes that should be credited for the reduced defects, then ideally the defect-reducing plans should be associated with a high precision, and among those plans, most of them should also have a high recall.

	Xtree	Shat	Oliv	Alves	Random	TimeLIME	LIME
Jedit	38	0	7	5	7	65	58
Camel1	80	0	8	7	7	90	73
Camel2	22	0	3	6	5	83	66
Log4j	23	1	10	8	12	81	70
Xalan	34	1	1	2	1	81	82
Ant	21	0	10	8	1	50	49
Velocity	19	0	2	3	13	88	85
Poi	45	1	4	18	4	61	76
Synapse	39	0	18	18	4	70	66
AVG	35.67	0.33	7.00	8.33	6.00	74.33	69.44
STD	17.99	0.47	5.01	5.52	4.03	12.79	10.69
Rank	2	4	3	3	3	1	1

TABLE 11: The precision rate (in percentage) of a plan measures how many changes proposed by the plan are found in the subsequent release. The rank is generated using the Scott-Knot test. A higher rank is better.

	Xtree	Shat	Oliv	Alves	Random	TimeLIME	LIME
Jedit	23	0	3	5	2	59	64
Camel1	29	0	2	3	3	59	52
Camel2	18	0	1	1	2	62	64
Log4j	5	1	8	6	6	58	72
Xalan	61	0	1	2	1	85	81
Ant	11	0	2	6	5	33	31
Velocity	27	0	1	3	2	73	77
Poi	30	0	0	6	2	60	56
Synapse	23	0	2	2	1	46	52
AVG	25.22	0.11	2.22	3.78	2.67	59.44	61.00
STD	14.88	0.31	2.20	1.87	1.63	13.85	14.46
Rank	2	4	3	3	3	1	1

TABLE 12: The recall rate (in percentage) of a plan measures out of all actual changes how many of them get proposed by the plan.

First, to evaluate if developers are capable of implementing changes proposed by our plans, we measure the **precision** rates of plans from different algorithms. The result from Table 11 shows that both LIME and TimeLIME planner obtain the highest scores. This is a supportive evidence indicating that developers, as seen in the subsequent release, were capable of implementing most of the changes proposed by our plans.

Secondly, to answer the question that whether or not developers may inadvertently make other changes while following plans proposed by planners, we measure the **recall** rates of plans. A low recall means that there exist more unanticipated changes. Therefore, it is more questionable that whether the plan or the unforeseen changes should take credit for the effect of defect reduction. As seen in Table. 12, TimeLIME and LIME still have the highest recall rates among all algorithms, which makes the performance of the planners more convincing since it is revealed here that when developers are making changes proposed by LIME or TimeLIME, they are less likely to deploy other changes that are not mentioned.

In summary, by mining the historical releases, the evaluational analysis here shows sufficient evidence that Time-LIME is of greater practicality compared to other algorithms. Furthermore, We also believe that more studies could be done to explore and expand the value of current evaluation process.

9 THREATS TO VALIDITY

Due to the complexity of the experiment designed in this case study, there are many factors that can threaten the validity of these results.

9.1 Learner Bias

This paper selects Random Forest classifier as the black-box classifier because prior research has shown that Random Forest classifier is ranked as one of the top models among all 32 classifiers used in defect prediction [58]. However, the preeminent predictive power of Random Forest classifier does not ensure that explanations derived from it are preeminent code refactoring plans as well. Other methods from the top rank may be more suitable in the problem of explanation generation while we haven't explored more.

9.2 Instrument Bias

Various approaches are proposed for explainable AI. Although LIME is one of the widely cited and well-known tools, it other tools might be suitable for solving SE problems, which can make solutions from LIME sub-optimal. Hence, to verify if adding in SE knowledge can always improve AI tools, we need to make a comprehensive exploration that includes more explanation generation methods.

9.3 Hyper-parameter Tuning

Past researches have shown how hyper-parameter optimization can boost the performance of a classifier used in defect prediction. Since in this paper we concentrate on the modification of the explainer instead of the learner, we used a simple grid search to find the optimal parameter setting. It can be possible that the current setting is sub-optimal and by using the actually optimal settings we might receive different experiment result.

10 RELATED WORK

Much research urges that interpretability should become an important factor in assessing analytical models in software engineering because software developers expect the model to provide understandable suggestions that can be actually achieved in real-world practice [64]–[66].

Recently at TSE'20, Jiarpakdee et al. modified LIME using hyper parameter optimization techniques, and assessed its performance in defect prediction via output stability [67]. The result has shown that explanations generated from their method are not only more stable among re-generations, but also understandable to software developers. The major difference is that:

- Jiarpakdee et al. assess the viability of applying modelagnostic techniques (such as LIME) in defect prediction whereas this paper assesses the practical effectiveness of LIME in re-organizing a project
- Jiarpakdee et al. explore possible means to improve the explanation generation procedure where as this paper explores methods to refine LIME's results into more actionable and effective plans for defect reduction.

11 FUTURE WORK

For future work, we need to take action to retire the above threats to validity.

11.1 More Learners

More black-box learners should be used in the experiment to construct a more comprehensive comparison. Although the limited sample amount of defect prediction datasets has ruled out many deep learning models such as Neural Network due to the overhead, there are still many other models, including but not limited to Random Subspace Sampling and Sequential Minimal Optimization, applicable for this experiment.

11.2 More Explainers

As described above, LIME is a representative member in the family of local surrogate interpretation models. Other local explanation generation methods that apply tree-structure extraction or association rule mining or so on should also be introduced in the discussion.

11.3 More Data

We would like to collect not only more SE projects of defection prediction data but also more releases of a single project. This can facilitate the further exploration on the accountability of our historical data analysis. According to the *K*-test, we validate plans on the more current release of the 3 releases. Because of that, we would prefer the files in the validation release is more similar to the proposed plans, no matter they have more or fewer bugs, so that our evaluation on the plans can be more accountable. Sometimes when the file in the validation release is not similar to the proposed plan, we wonder what would the file be like if there is another release with a more similar file. Could it be possible that more releases can provide us more accurate and robust evaluation conclusions?

11.4 Better Data

Our current data collection methods use all available training data. This means our model is learning from the track records of all developers. One potential issue with that approach is that the defect reduction performed by the development teams may not be the best. For example, say that in release 2, one team removed 5 defects with respect to release 1. Potentially, another team with more experience could have removed ten defects instead. Yet the above study learns from both teams even though one may be more effective than the other.

This is a fruitful area for useful future work. For example, in future work we could mine GitHub repositories looking for "tourists"; i.e. people that usually work in classes A,B,C,..., but then (occasionally) make changes to sections of the code that they are not so familiar with (class D,E, etc). Potentially, if we remove the "tourist" data and only learn from activities of more experienced developers, we might be able to build better plans using the better-refined training data.

11.5 More Measurements: Multi-objective Optimization

In this paper, we introduced *K*-test as a framework to conduction quality evaluation on changes proposed by different planners. However, although the current framework does provide us with insightful findings, we still believe that more measurements need to be brought in to construct a more comprehensive evaluation process. As shown in §8, planners make changes of different sizes, which makes it harder to examine their effectiveness since they are from different levels of precision: suppose 2 planners both made a change that overlaps with a defect-reducing action, the planner with a more precise/smaller change interval should probably be considered better than the other one. To address this problem, the future work could import another fitness score function that relates the precision to the effectiveness of the plan. That is to say, the task of planning defect reduction could be regarded as a multi-objective optimization problem, where a planner might have several goals (effectiveness, precision, feasibility, etc.) to chase after at the same time.

12 CONCLUSION

This paper has assessed the following *TimeLIME tactic* for generating defect reduction plans:

When reasoning about changes to a project, it is best to use changes seen in the historical record of that project.

Using this tactic, we find plans that:

- Are succinct: In terms of the average size of recommended plans. The TimeLIME generally generates smaller plans than the classical LIME and RandomWalk in every project. The plans are also equivalently succinct compared to other benchmark methods in this paper. Smaller plans are preferred to larger plan since the latter can be faster to apply.
- Better resemble developers' actions: In terms of the overlap between the proposed plans and the developer actions in the upcoming release, plans proposed by TimeLIME better match what developers actually do.
- Are better at reducing defects: In terms of the scaled weighted scores S_{scaled} that indicate the overall net gain received per project. TimeLIME gets the highest score among all planners in 8 out of 9 trials. (while the classical LIME wins in only 1 project).

Our results are a cautionary tale to the SE community. SE researchers need to be more careful about using off-theshelf AI tools, without first tuning them with SE knowledge. Specifically:

It is unwise to throw standard AI tools at SE problems without first considering how those tools might be customized for SE applications.

We say that since, as shown here, (a) such customization is not a complex thing to do and (b) the customized system can have dramatically better performance.

ACKNOWLEDGEMENTS

This work was partially funded by a research grant from the National Science Foundation (CCF #1703487) and the Laboratory for Analytical Sciences, North Carolina State University.

REFERENCES

- M. T. Ribeiro, S. Singh, and C. Guestrin, "" why should i trust you?" explaining the predictions of any classifier," in *Proceedings of the 22nd ACM* SIGKDD international conference on knowledge discovery and data mining, 2016, pp. 1135–1144.
- [2] A. Begel, J. Bosch, and M.-A. Storey, "Social networking meets software development: Perspectives from github, msdn, stack exchange, and topcoder,' IEEE Software, vol. 30, no. 1, pp. 52-66, 2013.
- D. Chen, K. Stolee, and T. Menzies, "Replication can improve prior results: [3] a github study of pull request acceptance," in 2019 IEEE/ACM 27th International Conference on Program Comprehension (ICPC). IEEE, 2019, pp. 179-190.
- [4]P. Devanbu, T. Zimmermann, and C. Bird, "Belief & evidence in empirical software engineering," in 2016 IEEE/ACM 38th International Conference on Software Engineering (ICSE). IEEE, 2016, pp. 108–119.
- C. Bird, N. Nagappan, B. Murphy, H. Gall, and P. Devanbu, "Don't touch my [5] code! examining the effects of ownership on software quality," in Proceedings of the 19th ACM SIGSOFT symposium and the 13th European conference on Foundations of software engineering, 2011, pp. 4–14.
- C. Passos, A. P. Braun, D. S. Cruzes, and M. Mendonca, "Analyzing the im-[6] pact of beliefs in software project practices," in 2011 International Symposium
- on Empirical Software Engineering and Measurement. IEEE, 2011, pp. 444–452. N. Shrikanth and T. Menzies, "Assessing practitioner beliefs about software [7] defect prediction," arXiv, pp. arXiv-1912, 2019.
- [8] Z. Wan, X. Xia, A. E. Hassan, D. Lo, J. Yin, and X. Yang, "Perceptions, expectations, and challenges in defect prediction," IEEE Transactions on Software Engineering, 2018.
- [9] T. L. Alves, C. Ypma, and J. Visser, "Deriving metric thresholds from benchmark data," in 2010 IEEE International Conference on Software Maintenance. IEEE, 2010, pp. 1–10. [10] R. Shatnawi, "A quantitative investigation of the acceptable risk levels of
- object-oriented metrics in open-source systems," IEEE TSE, 2010.
- [11] P. Oliveira, M. T. Valente, and F. P. Lima, "Extracting relative thresholds for source code metrics," in 2014 Software Evolution Week-IEEE Conference on Software Maintenance, Reengineering, and Reverse Engineering (CSMR-WCRE). IEEE, 2014, pp. 254-263.
- [12] S. Alexander, "refactoring.guru," http://regactoring.guru, Jan. 2014.
 [13] W. Fu, T. Menzies, and X. Shen, "Tuning for software analytics: Is it really necessary?" IST, 2016.
- [14] L. C. Briand, V. Brasili, and C. J. Hetmanski, "Developing interpretable models with optimized set reduction for identifying high-risk software components," IEEE Transactions on Software Engineering, vol. 19, no. 11, pp. 1028-1044, 1993.
- [15] F. Rahman, S. Khatri, E. T. Barr, and P. Devanbu, "Comparing static bug finders and statistical prediction," in ICSE. ACM, 2014.
- [16] Y. Shin and L. Williams, "Can traditional fault prediction models be used for vulnerability prediction?" EMSE, 2013. [Online]. Available: https://doi.org/10.1007/s10664-011-9190-8
- [17] T. Menzies, Z. Milton, B. Turhan, B. Cukic, Y. Jiang, and A. Bener, "Defect prediction from static code features: Current results, limitations, new ap-[18] T. Menzies, J. Greenwald, and A. Frank, "Data mining static code attributes
- to learn defect predictors," TSE, 2007.
- [19] C. Bird, N. Nagappan, H. Gall, B. Murphy, and P. Devanbu, "Putting it all together: Using socio-technical networks to predict failures," in ISSRE, 2009.
- [20] A. Arcuri and L. Briand, "A practical guide for using statistical tests to assess randomized algorithms in software engineering," in ICSE. IEEE, 2011.
- [21] D. Chen, W. Fu, R. Krishna, and T. Menzies, "Applications of psychological science for actionable analytics," *Foundations of Software Engineering*, 2018.
- [22] C. Rosen, B. Grawi, and E. Shihab, "Commit guru: Analytics and risk prediction of software commits," ser. ESEC/FSE 2015, 2015.
- [23] R. Krishna and T. Menzies, "Bellwethers: A baseline method for transfer learning," IEEE Transactions on Software Engineering, vol. 45, no. 11, pp. 1081-1105, 2018.
- [24] J. Nam, W. Fu, S. Kim, T. Menzies, and L. Tan, "Heterogeneous defect prediction," IEEE TSE, 2018.
- [25] B. Ghotra, S. McIntosh, and A. E. Hassan, "Revisiting the impact of classification techniques on the performance of defect prediction models," in 2015 37th ICSE.
- [26] A. Agrawal, W. Fu, and T. Menzies, "What is wrong with topic modeling? and how to fix it using search-based software engineering," IST, 2018.
- [27] A. Agrawal and T. Menzies, "Is better data better than better data miners?: on the benefits of tuning smote for defect prediction," in IST. ACM, 2018.
- [28] M. Fowler, Refactoring: improving the design of existing code. Addison-Wesley Professional, 2018.
- [29] W. C. Wake, Refactoring workbook. Addison-Wesley Professional, 2004.
- [30] T. Mens and T. Tourwé, "A survey of software refactoring," IEEE Transactions on software engineering, vol. 30, no. 2, pp. 126-139, 2004.
- [31] B. Du Bois, S. Demeyer, and J. Verelst, "Refactoring-improving coupling and cohesion of existing code," in 11th working conference on reverse engineering. IEEE, 2004, pp. 144–151. [32] M. Alshayeb, "Empirical investigation of refactoring effect on software
- quality," Information and software technology, vol. 51, no. 9, pp. 1319-1326, 2009.
- [33] B. Geppert, A. Mockus, and F. Robler, "Refactoring for changeability: A way to go?" in 11th IEEE International Software Metrics Symposium (METRICS'05). IEEE, 2005, pp. 10-pp.

IEEE TRANSACTIONS ON SOFTWARE ENGINEERING

- [34] R. Moser, A. Sillitti, P. Abrahamsson, and G. Succi, "Does refactoring improve reusability?" in *International Conference on Software Reuse*. Springer, 2006, pp. 287–297.
- [35] F. Simon, F. Steinbruckner, and C. Lewerentz, "Metrics based refactoring," in Proceedings Fifth European Conference on Software Maintenance and Reengineering. IEEE, 2001, pp. 30–38.
- [36] F. Dandashi, "A method for assessing the reusability of object-oriented code using a validated set of automated measurements," in *Proceedings of the 2002* ACM symposium on Applied computing, 2002, pp. 997–1003.
- [37] M. Bruntink and A. van Deursen, "An empirical study into class testability," Journal of systems and software, vol. 79, no. 9, pp. 1219–1232, 2006.
- [38] L. Tahvildari, "Quality-driven object-oriented re-engineering framework," in 20th IEEE International Conference on Software Maintenance, 2004. Proceedings. IEEE, 2004, pp. 479–483.
- [39] D. Wilking, U. F. Kahn, and S. Kowalewski, "An empirical evaluation of refactoring," e-Informatica Software Engineering Journal, vol. 1, no. 1, 2007.
- [40] B. Du Bois, S. Demeyer, and J. Verelst, "Does the" refactor to understand" reverse engineering pattern improve program comprehension?" in *Ninth European Conference on Software Maintenance and Reengineering*. IEEE, 2005, pp. 334–343.
 [41] R. Bender, "Quantitative risk assessment in epidemiological studies investi-
- [41] R. Bender, "Quantitative risk assessment in epidemiological studies investigating threshold effects," *Biometrical Journal: Journal of Mathematical Methods* in *Biosciences*, vol. 41, no. 3, pp. 305–319, 1999.
- [42] R. Krishna and T. Menzies, "From prediction to planning: Improving software quality with belltree," Empirical Software Engineering, 2020.
- [43] J. Han, H. Cheng, D. Xin, and X. Yan, "Frequent pattern mining: current status and future directions," *Data mining and knowledge discovery*, vol. 15, no. 1, pp. 55–86, 2007.
- [44] S. T. Mueller, R. R. Hoffman, W. Clancey, A. Emrey, and G. Klein, "Explanation in human-ai systems: A literature meta-review, synopsis of key ideas and publications, and bibliography for explainable ai," 2019.
- [45] T. Menzies, R. F. Cohen, S. Waugh, and S. Goss, "Applications of abduction: Testing very long qualitative simulations," *IEEE Trans. on Knowl. and Data Eng.*, vol. 14, no. 6, p. 1362–1375, Nov. 2002. [Online]. Available: https://doi.org/10.1109/TKDE.2002.1047773
 [46] M. T. Ribeiro, S. Singh, and C. Guestrin, "Anchors: High-precision model-
- [46] M. T. Ribeiro, S. Singh, and C. Guestrin, "Anchors: High-precision modelagnostic explanations," in *Thirty-Second AAAI Conference on Artificial Intelli*gence, 2018.
- [47] J. Hihn and T. Menzies, "Data mining methods and cost estimation models: Why is it so hard to infuse new ideas?" in Automated Software Engineering Workshop (ASEW), 2015 30th IEEE/ACM International Conference on. IEEE, 2015, pp. 5–9.
- [48] V. B. Kampenes, T. Dybå, J. E. Hannay, and D. I. K. Sjøberg, "A systematic review of effect size in software engineering experiments," *IST*, 2007.
- [49] R. Rosenthal, H. Cooper, and L. Hedges, "Parametric measures of effect size," The handbook of research synthesis, vol. 621, no. 2, pp. 231–244, 1994.
- [50] I. H. Witten, E. Frank, and M. A. Hall, Data Mining: Practical Machine Learning Tools and Techniques, 3rd ed. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 2011.
- [51] AAAI, "Aaai 1990 spring symposium series reports," AI Magazine, vol. 11, no. 3, p. 27, Sep. 1990. [Online]. Available: https://www.aaai.org/ojs/ index.php/aimagazine/article/view/848
- [52] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "Smote: synthetic minority over-sampling technique," *Journal of artificial intelligence research*, vol. 16, pp. 321–357, 2002.
- [53] A. Agrawal and T. Menzies, "Is "better data" better than "better data miners"? on the benefits of tuning smote for defect prediction," in *Proceedings of the 40th International Conference on Software Engineering*, ser. ICSE '18. New York, NY, USA: Association for Computing Machinery, 2018, p. 1050–1061. [Online]. Available: https://doi.org/10.1145/3180155.3180197
- [54] M. Jureczko and L. Madeyski, "Towards identifying software project clusters with regard to defect prediction," in *Proceedings of the 6th International Conference on Predictive Models in Software Engineering*, 2010, pp. 1–10.
- [55] N. Nagappan and T. Ball, "Static analysis tools as early indicators of prerelease defect density," in *Proceedings*. 27th International Conference on Software Engineering, 2005. ICSE 2005. IEEE, 2005, pp. 580–586.
- [56] M. Jureczko, "Significance of different software metrics in defect prediction," Software Engineering: An International Journal, vol. 1, no. 1, pp. 86–95, 2011.
- [57] L. Madeyski and M. Jureczko, "Which process metrics can significantly improve defect prediction models? an empirical study," *Software Quality Journal*, vol. 23, no. 3, pp. 393–422, 2015.
- [58] B. Ghotra, S. McIntosh, and A. E. Hassan, "Revisiting the impact of classification techniques on the performance of defect prediction models," in Proceedings of the 37th International Conference on Software Engineering-Volume 1. IEEE Press, 2015, pp. 789–800.
- [59] T. K. Ho, "Random decision forests," in Proceedings of 3rd international conference on document analysis and recognition, vol. 1. IEEE, 1995, pp. 278– 282.
- [60] M. Jureczko and D. Spinellis, "Using object-oriented design metrics to predict software defects," *Models and Methods of System Dependability. Oficyna Wydawnicza Politechniki Wrocławskiej*, pp. 69–81, 2010.
 [61] C. Couto, P. Pires, M. T. Valente, R. S. Bigonha, and N. Anquetil, "Predicting
- [61] C. Couto, P. Pires, M. T. Valente, R. S. Bigonha, and N. Anquetil, "Predicting software defects with causality tests," *Journal of Systems and Software*, vol. 93, pp. 24–41, 2014.
- [62] E. Shihab, Z. M. Jiang, W. M. Ibrahim, B. Adams, and A. E. Hassan, "Understanding the impact of code and process metrics on post-release defects: a case study on the eclipse project," in *Proceedings of the 2010 ACM-IEEE*

International Symposium on Empirical Software Engineering and Measurement, 2010, pp. 1–10.

- [63] T. M. Khoshgoftaar, E. B. Allen, R. Halstead, G. P. Trio, and R. M. Flass, "Using process history to predict software quality," *Computer*, vol. 31, no. 4, pp. 66–72, 1998.
- [64] T. Menzies and T. Zimmermann, "Software analytics: so what?" IEEE Software, vol. 30, no. 4, pp. 31–37, 2013.
- [65] C. Lewis, Z. Lin, C. Sadowski, X. Zhu, R. Ou, and E. J. Whitehead, "Does bug prediction support human developers? findings from a google case study," in 2013 35th International Conference on Software Engineering (ICSE). IEEE, 2013, pp. 372–381.
- [66] H. K. Dam, T. Tran, and A. Ghose, "Explainable software analytics," in Proceedings of the 40th International Conference on Software Engineering: New Ideas and Emerging Results, 2018, pp. 53–56.
- [67] J. Jiarpakdee, C. Tantithamthavorn, H. K. Dam, and J. Grundy, "An empirical study of model-agnostic techniques for defect prediction models," *IEEE Transactions on Software Engineering*, 2020.



Kewen Peng is a second year Ph.D. student in Computer Science at North Carolina State University. His research interests include using and refining artificial intelligence methods to solve problems in software engineering.



Tim Menzies (IEEE Fellow) is a Professor in CS at North Carolina State University. His research interests include software engineering (SE), data mining, artificial intelligence, search-based SE, and open access science. http://menzies.us