

Automatically Matching Bug Reports With Related App Reviews

Marlo Haering
University of Hamburg
Hamburg, Germany

haering@informatik.uni-hamburg.de

Christoph Stanik
University of Hamburg
Hamburg, Germany

stanik@informatik.uni-hamburg.de

Walid Maalej
University of Hamburg
Hamburg, Germany

maalej@informatik.uni-hamburg.de

Abstract—App stores allow users to give valuable feedback on apps, and developers to find this feedback and use it for the software evolution. However, finding user feedback that matches existing bug reports in issue trackers is challenging as users and developers often use a different language. In this work, we introduce *DeepMatcher*, an automatic approach using state-of-the-art deep learning methods to match problem reports in app reviews to bug reports in issue trackers. We evaluated *DeepMatcher* with four open-source apps quantitatively and qualitatively. On average, *DeepMatcher* achieved a hit ratio of 0.71 and a Mean Average Precision of 0.55. For 91 problem reports, *DeepMatcher* did not find any matching bug report. When manually analyzing these 91 problem reports and the issue trackers of the studied apps, we found that in 47 cases, users actually described a problem before developers discovered and documented it in the issue tracker. We discuss our findings and different use cases for *DeepMatcher*.

Index Terms—app store analytics, natural language processing, deep learning, mining software repositories, software evolution

I. INTRODUCTION

The app market is highly competitive and dynamic. *Google Play Store* and *Apple App Store* offer together more than ~ 4 million apps [46] to users. In this market, it is essential for app vendors to regularly release new versions to fix bugs and introduce new features [33], as unsatisfied users are likely to look for alternatives [8], [50]. User dissatisfaction can quickly lead to the fall of even popular apps [26]. It is thus indispensable to continuously monitor and understand the changing user needs and habits for a successful app evolution.

However, identifying and understanding user needs and encountered problems is challenging as users and developers work in different environments and have different goals in mind. On the one hand, software developers professionally report bugs in issue trackers to document and keep track of them, as illustrated in Figure 1. On the other hand, users voluntarily provide feedback on apps in, e.g., app reviews as shown in Figure 2 – using a different, often non-technical, and potentially imprecise language. Consequently, seriously considering and using app reviews in software development

and evolution processes can become time-consuming and error-prone.

App vendors can regularly receive a large number of user feedback via various channels, including app stores or social media [10], [37]. Manually filtering and processing such feedback is challenging. In recent years, research developed approaches for filtering feedback, e.g., by automatically identifying relevant user feedback [4] like bug reports [44] and feature requests [16], or by clustering the feedback [47] to understand how many users address similar topics [49]. While these approaches are helpful to cope with and aggregate large amounts of user feedback, the gap between what happens in the issue tracker and what happens online in the user space remains unfilled. For instance, developers remain unable to easily track whether an issue reported in an app review is already filed as a bug report in the issue tracker; or to quickly find a related bug they thought is already resolved. Additionally, user feedback items often lack information that is relevant for developers, such as steps to reproduce or versions affected [31], [53].

To address this gap, we introduce *DeepMatcher*, which is, to the best of our knowledge, the first approach that matches official and technically-written bug reports with informal, colloquially-written app reviews. *DeepMatcher* first filters app reviews into **problem reports** using the classification approach by Stanik et al. [44]. Subsequently, our approach matches the problem reports with **bug reports** in issue trackers using deep learning techniques. We use the state-of-the-art, context-sensitive text embedding method DistilBERT [41] to transform the problem report and bug report texts into the same vector space. Given their vector embeddings, we then use cosine similarity as a distance metric to identify matches.

For 200 randomly sampled problem reports submitted by users of four Google apps, *DeepMatcher* identified 167 matching bug reports when configured to show three suggestions per problem report. In about 91 cases, *DeepMatcher* did not find any matches. We manually searched for these 91 cases in the issue trackers to check whether there are indeed no matching bug reports. We found that in 47 cases, developers would have benefited from *DeepMatcher*, as no corresponding bug reports were filed. We also qualitatively analyzed the context-sensitive text embeddings, which identified recurring bug reports and cases in which users reported problems before

-  **Deleting old messages crashes app - Video**
 #9920 opened 10 hours ago by condemnedmeat

-  **Message reactions from another user show up in my reaction menu on the same message**
 #9918 opened 17 hours ago by kyleedwardsny

-  **Member does not support new groups still when using the same beta version**
 #9917 opened 17 hours ago by Edu4rdSHL

-  **Reaction notification not dropped on phone once read on desktop.**
 #9909 opened 5 days ago by iago-lito

Fig. 1. List of bug reports from the issue tracker of the app Signal Messenger.

developers documented them. We found that our approach can detect semantically similar texts such as “draining vs. consuming battery” and “download vs. save PDF”, filling the gap between users’ and developers’ language. Our qualitative analysis further revealed cases of recurring and duplicated bug reports. We share our replication package¹ for reproducibility.

The remainder of the paper is structured as follows. First, we introduce *DeepMatcher* in Section II explaining our design rationales. Section III introduces our evaluation setting, including the research questions, data, and process. Section IV presents our quantitative and qualitative evaluation results. Then, we discuss how developers can use and modify *DeepMatcher* to detect bugs earlier and enrich existing issue descriptions with information extracted from user feedback in Section V. Finally, we discuss the threats to validity in Section VI, related work in Section VII, and conclude the paper in Section VIII.

II. APPROACH

Figure 3 shows an overview of *DeepMatcher*’s technical approach. The input of *DeepMatcher* is a problem report (an app review describing a problem with an app) and a bug report summary. In Section II-A, we discuss how we automatically identified problem reports from the review. Section II-B describes the text embedding creation process shown in the middle part of the figure. This represents the transformation of textual data into numeric values, which we then use to calculate a similarity value as explained in Section II-C.

A. Automatic Problem Reports Classification

Challenges. One of the major problems when working with user feedback is the vast amount that software developers receive. Particularly in app stores, Pagano and Maalej [37] showed that developers of popular apps receive about 4,000 app reviews daily. When considering Twitter as an alternative feedback source, Twitter accounts of popular software vendors receive about 31,000 tweets daily [10]. Besides the amount, the quality of the written feedback differs. Most of the received app reviews simply praise, e.g., “I like this app” or

¹<https://mast.informatik.uni-hamburg.de/replication-packages/>



Fig. 2. Example problem report for Nextcloud answered by the app developer.

dispraise, e.g. “I hate this app!” [37]. However, developers are particularly interested in the user experience, feature requests, and problem reports [11], [27], [47]. Our approach uses automatically classified problem reports from app reviews and subsequently matches them to bug reports in issue trackers.

Approach and Rationale. We applied a four-step process to filter relevant app reviews. First, we removed all user feedback containing less than ten words as previous research has shown that such feedback is most likely praise or spam [37] and does not contain helpful information [42]. Second, we downloaded the replication package of Stanik et al. [44], and applied the *bug report*, *feature request*, and *irrelevant* classification approach to also filter the user feedback for bug reports. The classification reduced the initial number of app reviews to 9,132 problem reports. Fourth, to check the reliability of the classification, we randomly sampled and manually analyzed automatically classified app reviews for each of the four studied apps for manual analysis. Two coders manually checked if the classified problem reports were correctly classified. In case of disagreement, we did not include the app review but sampled a new one. We repeated this step until we had 50 verified problem reports for each app, which is 200 in total.

B. Text Representation with Word Embeddings

Challenges. We further convert the text into a numerical presentation for further interpretation. In natural language processing, practitioners usually transfer texts into vectors by applying techniques including bag-of-words, tf-idf [29], or fastText [20]. When representing text in a vector space, we can perform calculations, such as comparing text similarity, which becomes essential in a later step for identifying matches. Selecting the right word embedding technique is crucial, as it decides how precisely the vectors represent the text. We face two major challenges. First, users and developers usually use different vocabularies. User feedback is more prone to spelling mistakes, often contain emoticons. Moreover, users write mostly in an informal, colloquial, and non-technical way. Second, bug reports are usually written in a more formal way, e.g., following templates, containing metadata, and may provide technical information like stack traces [53].

Approach and Rationale. Both data sources consist of different text components. While user feedback consists of a single text body, bug reports have a summary and a detailed

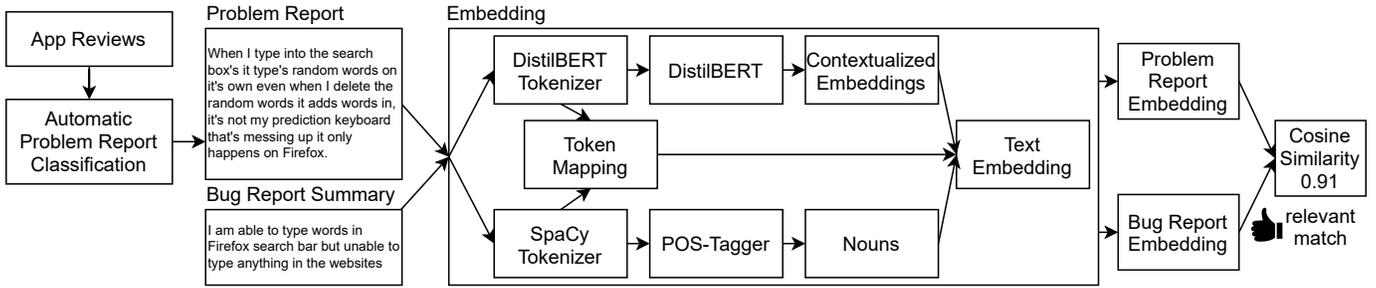


Fig. 3. Overview of the *DeepMatcher* approach.

description. The description may contain a long explanation, including steps to reproduce, stack traces, and error logs. We determined which text components of the bug report to include for the calculation of the text embedding. Previous research showed that the detailed bug report description contains noise for information retrieval tasks [51]. In particular, it contains technical details that users usually omit in their user feedback [31]. Further, research shows that the summary already contains the essential content of the long description [22], [25], [48]. Therefore, we calculated the word embeddings only based on the bug report’s summary.

Regarding the word embedding technique, we chose DistilBERT [41], a light-weight version of BERT [6] that is trained with a fewer number of parameters but has a similar generalization potential. Alternative techniques would be, e.g., BERT, XLNet, or RoBERTa. But as DistilBERT requires significantly fewer hardware resources and training time, it is more applicable for various development teams. Our technique first tokenizes the input text and then calculates vectors for each token. Compared to other text representations like bag-of-words or tf-idf, these vectors are contextualized; they consider the context of the surrounding words. For example, the two sentences “I love apples” and “I love Apple macbooks” contain the token “apple”. Contextualized embeddings take into account that the token’s semantics differs in these two sentences. In our approach, DistilBERT creates a 768-dimensional contextualized embedding for each token.

We calculated the document embedding from the individual token embeddings. To reduce the weight of frequent but unimportant words such as “I”, “have”, or “to”, previous research in text mining suggests removing stopwords [38], [45], [47]. In our approach, we went one step further and only included embeddings of nouns, which we can automatically detect with a part-of-speech (POS) tagger. We carefully decided to remove other parts of speech like the verb tokens as first trials showed that including frequent verbs like “crash”, “freeze”, and “hangs” heavily biased our results toward these terms. For example, *DeepMatcher* would match “The app crashes when I open a new tab” (problem report) with “Firefox crashes on the home screen” (bug report) because the verb “crash” puts the vectors of both texts closer together. Based on this design decision, *DeepMatcher* weights essential words, i.e., nouns that describe components or features higher, while the contex-

tualized token embeddings still contain information about the surrounding context, e.g., the verbs. As a result, *DeepMatcher*, e.g., emphasizes the nouns “new tab” and “home screen” in the previous example and, therefore, would not consider the bug and problem report as a potential match. Another positive side-effect of the surrounding context is that it helps to deal with misspelled words as their surrounding context is usually similar to the correct word’s context. Therefore DistilBERT calculates similar embeddings for them.

The automatic noun detection of the input texts is part of *DeepMatcher* and uses SpaCy’s tokenizer and POS-tagger [15]. As SpaCy’s tokenizer and the DistilBERT’s tokenizer split the input text into different token sequences, we mapped the two sequences to each other by aligning them using pytokenizations. For calculating the embedding for the full text of the problem report or bug report’s summary, we added all noun word vectors of the text and averaged them. Alternatively, we could have summed up the noun word vectors but decided to average them as the cosine similarity function depends on the vector angles and not on their lengths. Therefore, the choice of summing or averaging would not influence the cosine similarity score in our approach.

C. Identifying relevant Bug Reports for a Problem Report

Challenges. Given the numerical representation of the problem report and the bug report, *DeepMatcher* finally requires a method to decide whether a bug report is relevant for a problem report or not. The main challenge in this task is calculating matching problem reports and bug reports with *short text similarity* [40]. Besides semantic features, research tried text similarity approaches like simple substring comparisons [17], lexical overlap [18], or edit distances [34]. Kenter and de Rijke [21] state that these approaches can work in some simple cases but are prone to mistakes.

Approach and Rationale. We considered two options for this task. One option is to model this task as a binary classification problem using the two classes, “relevant” or “not relevant”. However, this approach would require a large labeled dataset to train a classifier for this task, which is expensive and time-consuming [5]. Therefore, we chose the second option, which models this task as an information retrieval task. Given a problem report as a query, we designed *DeepMatcher* to return a ranked list of suggested relevant bug reports. We chose a

distance function to measure the similarity between the two text embeddings and further rank the bug report summaries in decreasing order.

Two popular similarity measures for text embeddings are the euclidian similarity, and cosine similarity [9], [43]. The euclidian distance can increase depending on the number of dimensions. In contrast, the cosine similarity measures the angle of the two text vectors and is independent of their magnitude. The benefit is that it results in a similarity value of 1 if the two vectors have a zero-degree angle. A non-similarity occurs when the vectors have a 90-degree angle to each other. Previous research [1]–[3], [9], showed that cosine similarity performs equally or outperforms other similarity measures for dense text embedding vectors, which is why we also used it for *DeepMatcher*.

III. EMPIRICAL EVALUATION

A. Research Questions

Apps usually receive user feedback as app reviews, which may contain the user’s opinion and experience with the software. Our study focuses on the app review category *problem reports*, which is about users describing a faulty app behavior. Figure 2 shows an example of a problem report. *Bug reports* are issues in an issue tracker, complying with a certain template and contain information including summary, body, app version, timestamp, and steps to reproduce. Figure 1 shows a list of four bug report summaries of the Signal Messenger app in GitHub.

Our evaluation focuses on the following research questions:

RQ1 How accurate can *DeepMatcher* match app reviews with bug reports?

In this research question, we analyze if we can identify bug reports in issue tracker systems for which users wrote app reviews. For example, a developer filed the following bug report: “*I am able to type words in Firefox search bar but unable to type anything in the websites*”. Can we find app reviews that describe the same issue, like the following app review? “*When I type into the search box’s it type’s random words on it’s own even when I delete the random words it adds words in, it’s not my prediction keyboard that’s messing up it only happens on Firefox.*”

RQ2 What can we learn from *DeepMatcher*’s relevant and irrelevant matches?

To answer this question, we checked a sample of relevant and irrelevant matches. We analyzed cases in which contextual embeddings identify words with similar meanings, the language gap between developers and users, recurring bug reports, and a potential chronological dependency between problem reports and bug reports. We highlight our findings and explain them with examples from our dataset.

B. Evaluation Data

For creating the evaluation data, we first collected app reviews and issues of four diverse apps. We selected the apps Firefox (browser), VLC (media player), Signal (messenger),

TABLE I
OVERVIEW OF THE EVALUATION DATA.

App name	Bug reports		App reviews		
	Time period	Count	Time period	Count	Problem reports [44]
Firefox Browser	01/2011 08/2019	29,941	09/2018 07/2020	5,706	3,314
VLC Media Player	05/2012 07/2020	553	09/2018 07/2020	5,026	2,988
Signal Messenger	12/2011 08/2020	7,768	09/2018 08/2020	10,000	2,583
Nextcloud	06/2016 08/2020	2,462	06/2016 08/2020	774	247
Total		40,724		21,506	9,132

and NextCloud (cloud storage) to cover different app domains and usage scenarios in our analysis. As Pagano and Maalej showed [37], most app reviews rather represent noise for software developers as they only contain praise like “I like this app” or insults like “this app is trash”. Therefore, we applied the bug report classification approach from Stanik et al. [44] to identify problem reports. We chose this classification approach, as it uses state-of-the-art approaches, achieves high F1-scores of 79% for English app reviews, and we could include the replication package in our pipeline without major modifications. Eventually, we created a random sample from the collected data that we then used in the evaluation.

As our study is concerned with matching problem reports found in app reviews with bug reports documents in issue trackers, we collected both—problem reports and bug reports. In our study, we decided to evaluate our approach against four popular open-source Android apps, which stretch over different app categories. We cover the categories browsing (Firefox), media player for audio, video, and streaming (VLC), a cloud storage client (Nextcloud), and a messaging app (Signal). As these apps use different issue tracker systems to document bug reports, we developed crawlers for Bugzilla (Firefox), Trac (VLC), and GitHub (Nextcloud and Signal). For each app, we collected all bug reports from the issue tracker systems. As a requirement for our analysis, each bug report contains at least an ID, summary, and status (e.g., open and resolved), as well as the creation date. Additionally, we also collected the remaining data fields provided by the issue tracker systems, such as issue descriptions and comments. A complete list of the collected data fields is documented in our replication package.

We then collected up to 10,000 app reviews of the corresponding apps following Google’s default sort order “by helpfulness”. Sorting by helpfulness helped us to not only considering the most recent app reviews (sort by date) but also emphasized the app reviews that other users deemed helpful. For Nextcloud, we could not collect more than 774 app reviews as it seems that from their total 5,900 reviews in the Google Play Store, 5,126 reviews only contain a star rating without any text. Our app review dataset covers a time frame of two to four years. In total, we were able to collect 21,506 app reviews from the Google Play Store. After applying the problem report

classifier [44], we could reduce the number of app reviews to 9,132 problem reports.

Table I summarizes our study data. The table reveals that while the time range of bug reports covers at least four years, Firefox has the highest number of bug reports filed from January 2011 to August 2019. In total, we collected 40,724 bug reports, of which 29,941 belong to Firefox. We focused on bug reports but ignored other issues like feature or enhancement requests by filtering the issues that developers labeled as such in the issue tracker systems.

C. Evaluation Method

We evaluated *DeepMatcher* with respect to quantitative and qualitative aspects. Starting from a set of manually verified problem reports, *DeepMatcher* suggested three bug reports for each. We evaluated how accurately *DeepMatcher* finds matching bug reports based on their summaries for a given problem report. We conducted a manual coding task, which consisted of two steps.

In the first step, we classified the app reviews using an existing approach [44] into problem reports to remove irrelevant feedback such as praise and dispraise (F1-score of 0.79 for English app reviews). Then, we randomly sampled 50 problem reports per app and manually verified whether the classified app reviews are problem reports. Two coders independently annotated the classification results according to a coding guide from previous research [27]. We randomly sampled new problem reports until we reached 50 verified problem reports per app, which made 200 in total.

In the second step, we used *DeepMatcher* to calculate three suggestions of potentially matching bug reports for each of the 200 problem reports. Again, two coders independently read each problem report and the three suggested bug reports. For each matching, the coders annotated whether the match is relevant or irrelevant. We consider the match relevant if the problem report and the bug report describe the same app feature (e.g., watch video) and behavior (e.g., crashes in full screen). For example, for the problem report: “*Latest update started consuming over 80% battery. Had to uninstall to even charge the phone!*” *DeepMatcher* suggested the relevant bug report match “*Only happening with latest version, But keep getting FFbeta draining battery too fast*”. We documented the inter-coder agreement and resolved disagreements by having the two coders discussing each. We report further analysis results based on the resolved annotations.

To answer RQ1, we calculated *DeepMatcher*’s performance. We report the number of relevant/irrelevant matches found per app and the mean average precision (MAP) [52]. It describes the average precision *AveP* for each problem report p and its suggestions and then calculates the mean over all problem reports P :

$$MAP = \frac{\sum_{p=1}^P AveP(p)}{P}$$

This is a conservative evaluation metric because it assumes that we have at least three relevant bug reports per problem

report. If this is not the case, even a perfect tool cannot achieve the highest average precision [30]. However, in our setting, the actual number of relevant bug reports is unknown. Therefore, we additionally report on the hit ratio, which describes the share of problem reports for which *DeepMatcher* has suggested at least one relevant match. For the irrelevant matches, we further tried to manually find relevant bug reports in issue trackers. We further analyzed *DeepMatcher*’s similarity score to identify a possible threshold, which users can use for the relevance assessment.

To answer RQ2, we conducted a qualitative analysis of the data. For each app, we analyzed the language of app reviews and bug reports by counting the nouns used in both datasets in relation to the nouns used overall. We highlight the strength of contextual word embeddings and show how *DeepMatcher* matches different words with similar semantic meaning. We further analyze the cases in which developers report a bug report after a user submitted a related problem report in the app store.

IV. EVALUATION RESULTS

This section reports the results of our evaluation study. We analyze *DeepMatcher*’s cosine similarity values to understand if we could use a certain similarity score threshold to identify matching problem reports and bug reports. Further, we report on our qualitative analysis and describe relevant and irrelevant suggestions to find potential ways to improve automatic matching approaches.

A. Matching Problem Reports with Bug Reports (RQ1)

As introduced earlier, we sampled 50 problem reports per app (200 in total) and applied *DeepMatcher* to suggest matching bug reports. In the first step, *DeepMatcher* suggested one matching bug report per problem report. Then, we changed that parameter and let *DeepMatcher* suggest two matching bug reports. Finally, *DeepMatcher* suggested three matching bug reports per problem report, which led to 600 suggestions. Since *DeepMatcher* suggests bug reports based on the highest cosine similarity, it added one additional suggestion per step while keeping the previous ones. This way, we could evaluate *DeepMatcher*’s performance based on this parameter (number of suggestions). Two authors independently annotated each of the 600 bug report suggestions as either a relevant or irrelevant match.

Table II summarizes the overall result of the peer-coding evaluation. The table shows that the inter-coder agreement for the whole dataset (3 suggested bug reports per problem report) is ≥ 0.88 . From the 600 matching bug report suggestions, the two coders identified 167 developer relevant matching suggestions. These 167 suggestions occurred in 109 problem reports with the parameter *number of suggestions* set to three. Multiple relevant matches occurred either for generic problem reports like “the app crashes often” or for similar, recurring, or duplicated bug reports in the issue tracker.

Suggestions Without Relevant Matches. For 91 problem reports, *DeepMatcher* could not find a relevant match within

TABLE II

RESULTS OF THE MANUAL CODING FOR 4 OPEN SOURCE APPS, EACH WITH 50 APP REVIEWS. LEGEND: MEAN AVERAGE PRECISION (MAP), NUMBER OF SUGGESTED BUG REPORTS (#).

App	1 Suggestion			2 Suggestions			3 Suggestions				Coder Agreement
	#	MAP	Hit Ratio	#	MAP	Hit Ratio	#	MAP	Hit Ratio	# Relevant Matches	
Firefox	50	0.50	0.50	100	0.54	0.58	150	0.58	0.74	38	0.93
VLC	50	0.32	0.32	100	0.38	0.44	150	0.40	0.51	26	0.91
Signal	50	0.38	0.38	100	0.47	0.57	150	0.50	0.68	45	0.89
Nextcloud	50	0.62	0.62	100	0.73	0.84	150	0.73	0.89	58	0.88
Total	200	∅ 0.45	∅ 0.46	400	∅ 0.53	∅ 0.61	600	∅ 0.55	∅ 0.71	167	∅ 0.90

the three suggestions. The reason for this is twofold: either no relevant bug report actually exists in the issue tracker system, or *DeepMatcher* missed relevant matches. To understand why *DeepMatcher* did not identify any matches for 91 problem reports, we manually searched the issue tracker systems by building a query using different keyword combinations from the problem reports. For example, Table III shows a problem report of VLC for which *DeepMatcher* could not find a relevant matching bug report. However, in our manual check, we found the bug report “When the device’s UI language is RTL, no controls are shown in the notification card”, which the two coders consider a relevant match. For 47 problem reports, we could not find any relevant match in the issue tracker system, while *DeepMatcher* missed potentially relevant matches in 44 cases. Consequently, *DeepMatcher* identified 47 problem reports that were undocumented in the issue trackers. This can help developers create new bug reports.

Average Mean Precision and Hit Ratio. We calculated the Mean Average Precision (MAP) and the hit ratio of our manual annotated data for all three parameters (one suggestion, two suggestions, and three suggestions). The MAP is a conservative score, which assumes that each problem report has at least as many relevant bug reports in the issue tracker as the parameter states. For example, if we set the parameter for the number of suggested bug reports to three, the MAP score assumes that at least three relevant matching bug reports exist. In case the problem report has less than three relevant bug reports, the average precision for that problem report cannot get the maximum value of one [30]. For our calculation, we excluded the problem reports for which we could not find a relevant bug report manually. The hit ratio, on the other hand, is the number of problem reports for which *DeepMatcher* found at least one relevant match divided by the number of all problem reports.

Table II shows the MAP and the hit ratio scores for each parameter setting. Increasing the parameter from one to two shows that the MAP score increases by 8%, while the hit ratio increases by 15%, which means that we increase the chance of finding a relevant match to 61%. When further increasing the parameter to three, we observe that the probability of having at least one relevant match increases to 71%, however as the MAP score reveals, developers might have to consider more irrelevant matches. We found that for Nextcloud, *DeepMatcher*

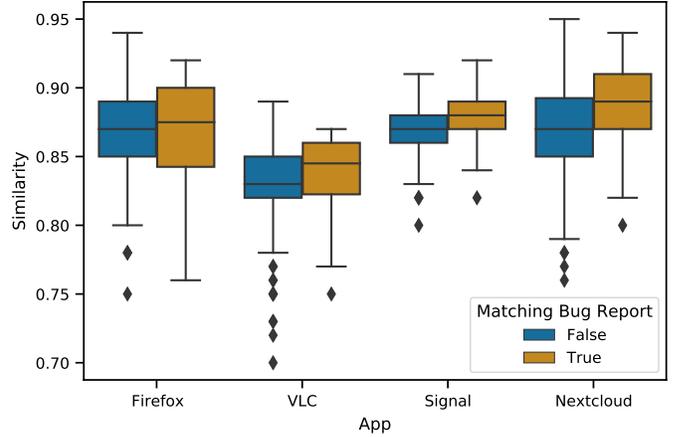


Fig. 4. Similarity values for relevant and irrelevant matches per app.

achieved the highest Mean Average Precision (0.73) and Hit Ratio (0.89). In contrast, VLC achieved the lowest scores with a MAP of 0.40 and a hit ratio of 0.51. Averaged over all apps, *DeepMatcher* achieved a mean average precision of 0.55 and a hit ratio of 0.71.

Cosine Similarity Analysis. We analyzed the cosine similarity values of relevant bug report suggestions and the irrelevant bug report suggestions. Figure 4 shows the cosine similarity values for the manual labeled suggestions for each app. It shows that the medians of the similarity scores of relevant bug report matches are higher than the irrelevant matches. However, the similarity scores vary between their min and max values by up to ~ 0.15 . All similarity scores are overall high (≥ 0.5) as all texts are in the technical domain.

We found that VLC has the lowest cosine similarity score compared to the other apps, which is also the app for which *DeepMatcher* found the fewest relevant bug report matches (26 matches). The lower cosine similarity indicates a higher language gap between VLC problem reports and bug reports. To further analyze this indication, we calculated the overlap of nouns used in problem reports with the nouns used in bug report summaries. We only checked the noun overlap as this is the part-of-speech category *DeepMatcher* uses to generate matches. For each app, we calculated the ratio between the number of nouns used in problem reports and bug reports, and the number of nouns used overall. The apps’ ratios are:

Firefox 19%, VLC 11%, Signal 24%, and Nextcloud 25%. The noun overlap calculation strengthens our assumption that the language between the VLC problem reports and the bug report summaries diverge more than the other apps, which negatively affects *DeepMatcher*'s automatic matching approach.

B. Qualitative Analysis of *DeepMatcher*'s Relevant, Wrong, and Missed Suggestions (RQ2)

We summarize and describe qualitative insights to learn about *DeepMatcher*'s relevant and irrelevant suggestions. Table III provides examples of problem reports, *DeepMatcher* suggested bug report summaries, and our coding of whether we think that there is a relevant match for developers. In the table, we selected one problem report per app and searched for cases that highlight some of our insights, like recurring bug reports for Signal, or a problem report submitted long before a bug report was filed in the Nextcloud app. In the following, we discuss our findings.

Strength of Contextual Embeddings. One strength of our approach is to learn the context of words (which words belong together). Other approaches like bag-of-words or tf-idf do not consider the context of words and, therefore, fall short in representing a deeper understanding of the language. The following two examples illustrate the strength of word context. *DeepMatcher* suggested matches that included the phrases “automatic synchronization” and “auto upload” in Nextcloud bug reports, as well as “download pdf” and “save pdf” for the Firefox browser. One full example is shown in Table III. The Firefox problem report discusses a “consuming battery” problem that happens since the “latest update”. The relevant matching bug report states that the “battery draining” becomes a problem in the “latest version”. It shows that the contextual embeddings of the noun tokens, e.g., “synchronization” and “upload” reach a high text similarity score as they are considered closely related.

Language Gap Leads to Fewer Relevant Matches. During the manual coding task, we noticed that the phrasings in VLC's bug reports often contain technical terms, for example: “Freeze entire android OS when playing a video. libvlc stream: unknown box type cTIM (incompletely loaded)”. However, users are typically not part of the development team and do not include technical words like specific library names used by the developers. Our previously reported plot of the cosine similarity values in Figure 4 quantitatively indicated that there might be a language gap as the text similarity scores between problem reports and bug reports were the lowest for VLC. We then performed a noun overlap analysis, which strengthened the indicator for the language gap as VLC has the lowest noun overlap with 11%. Eventually, we looked into the problem reports, bug reports, the Google Play Store, and the issue trackers.

We found that the developers of Nextcloud sometimes reply to problem reports in the Google Play Store and ask the users to also file a bug report in the issue tracker systems. We do not know how many users are actually going to the issue tracker

to report a bug. But this could also explain why Nextcloud has the highest cosine similarity score and highest noun overlap (25%). Consequently, *DeepMatcher* is more accurate if bug report summaries contain non-technical phrases, as users rarely use technical terms.

Sometimes users do not understand the app features. The following example from a Signal problem report shows a user confusing a feature with a bug: “Works well but gives me 2 check marks immediately after sending my text. I know the receivers are not reading the texts so fast. Why 2 checkmarks?” The two checkmarks in Signal are shown as soon as the addressed user successfully received the message. Signal has an optional feature that changes the color of the two checkmarks to blue if the recipient reads the message.

Recurring Bug Reports. Table III shows an example of recurring bug reports. The problem report of the Signal app states that the user did not receive notifications of incoming messages. We considered all three matching bug report suggestions of *DeepMatcher* as relevant, as they state the same problem. The interesting insight in this example is that with *DeepMatcher* we were able to identify a recurring bug report, as the first one was filed in September 2015, the second in January 2016, the third in April 2016. The problem report of the user happened in October 2017, more than one year after the last suggested match. In Section V-C, we discuss how *DeepMatcher* can help systematically find such cases.

Date Case Analysis. Regarding the date analysis, we found that in 35 of 167 relevant matches, developers reported the bug reports after the users submitted the corresponding review in the Google App Store. The time differences of the 35 cases in which the problem report submission happened before the bug report, is 490 days later, on average. In the following, we illustrate three examples.

Table III shows one problem report for Nextcloud, submitted in October 2017, while the matching bug report was filed in July 2018. Another user submitted the following problem report on the Nextcloud app: “Autoupload not working, android 7, otherwise all seems good. Happy with app and will increase stars to 5 when auto upload is working.” *DeepMatcher* identified the matching bug report “Android auto upload doesn't do anything” that was created 29 days after the problem report. In the last example, a developer documented a matching bug report 546 days after the corresponding problem report for the Signal app. Both the user and the developer address the in-app camera feature: “Newest update changes camera to add features, but drastically reduces quality of photos. Now it seems like the app just takes a screenshot of the viewfinder, rather than taking a photo and gaining from software post-processing on my phone. [...]”. The bug report stated: “In-app camera shows different images for preview and captured”.

V. DISCUSSION

This section discusses potential use cases of *DeepMatcher* to support developers in their software evolution process.

TABLE III
EXAMPLE PROBLEM REPORTS FROM APP REVIEWS AND *DeepMatcher*'s SUGGESTED MATCHING BUG REPORTS. THE RELEVANT COLUMN SHOWS WHETHER THE TWO CODERS ANNOTATED THE SUGGESTIONS AS RELEVANT FOR DEVELOPERS.

Problem Report	Suggested Bug Report Summary	Relevant
App: Firefox Date: 2020-04-21 Report: Latest update started consuming over 80% battery. Had to uninstall to even charge the phone!	Date: 2018-01-04 Report: Only happening with latest version, But keep getting FFbeta draining battery too fast	yes
	Date: 2017-11-20 Report: The topbar on android phone becomes white, which makes the time and battery life invisible.	no
	Date: 2016-12-13 Report: "Offline version" snackbar is displayed when device is very low on power and in battery saving mode	no
App: Signal Date: 2017-10-09 Report: it is a good app. i am mostly satisfied with it but sometimes, the notifications would not work; so, I would not know that someone messaged me until I open the app. it might have been fixed because it hasn't been happening in the last month or so. Would recommend.	Date: 2015-09-13 After update: no notification sent with TextSecure message. I have to open the app to see if there's something new	yes
	Date: 2016-01-17 Report: No notifications show up until the app is manually open	yes
	Date: 2016-04-17 Report: Not getting notification in real time unless I open the app	yes
App: VLC Date: 2020-05-17 Report: So many bugs... Plays in background, but no controls in notifications. When you tap the app to bring up the controls, the video is a still screen. Navigating is a pain. Resuming forgets my place constantly. Basically unusable	Date: 2019-04-13 Report: android navigation bar, shown after a click, shifts and resizes full-screen video	no
	Date: 2018-09-27 Report: Play/pause button icon is not shifting while pausing the audio on notification area	no
	Date: 2013-09-16 Report: [Android] On video playing the navigation bar is not hidden on some tablets	no
App: Nextcloud Date: 2017-10-09 Report: I have a nextcloud server and the way I access my server is via OpenVPN. The problem now is the nextcloud native app doesn't work through vpn. It is an odd behavior. I highly recommend to use owncloud app instead.	Date: 2016-07-09 Report: nextcloud android client can't login but android webdev clients do	no
	Date: 2018-07-20 Report: AutoUpload stuck on "Waiting for Wifi" when using VPN	yes
	Date: 2020-06-18 Report: SecurityException in OCFileListAdapter: uid 10410 cannot get user data for accounts of type: nextcloud	no

A. Detecting Bugs Earlier

It is essential for app developers to address users' problems as their dissatisfaction may lead to the fall of previously successful apps [26], [50]. One way to cope with user satisfaction is to quickly fix frustrating bugs, which may cause users to switch to a competitor and submit negative reviews. However, bugs may occur for different reasons. Some bugs affect only a few users with specific hardware or software versions, while others affect a large user group. Further, not all bugs are immediately known to developers, particularly, non-crashing bugs, which are hard to discover in automated testing and quality assurance [31]. Our results show that some users submit problem reports in the Google Play Store months before developers document them as bug reports in the project's issue tracker. When considering additional feedback channels such as social media and other stores, this might get even worse.

Our qualitative analysis of bug reports shows that these earlier submitted problem reports contain valuable information

for app developers, such as the affected hardware. Therefore, we emphasize that developers should continuously monitor user feedback in app stores to discover problems early and document them as bug reports in their issue trackers [32]. For this purpose, developers can first apply the automatic problem report classification of app reviews and subsequently use *DeepMatcher* to find existing matching bug reports. In case *DeepMatcher* does not find matching bug reports, we suggest that developers should consider the problem report as an unknown bug. However, to avoid the creation of duplicate bugs, we further suggest checking the issue tracker beforehand. Mezouar et al. [7] suggest a similar recommendation for developers when considering tweets instead of app reviews. They show that developers can identify bugs 8.4 days earlier for Firefox and 7.6 days earlier for Chrome (on average).

We envision different ways to suggest new bugs to developers. First, we could build a system that shows newly discovered bugs to developers. From that system, developers can decide to file a new issue in the issue tracker, delay, or reject it. Alternatively, a bot can, e.g., file a new issue in the

issue tracker systems automatically [31]. For the latter, future research could develop, e.g., approaches could prepare certain text artifacts, including steps to reproduce, meaningful issue description, or context information in a template for creating a new issue.

Furthermore, *DeepMatcher*'s application is not limited to user feedback in the form of app reviews. Our approach can generally process user feedback on various software, which developers receive via different channels, including app stores, social media platforms like Twitter and Facebook, or user support sites. *DeepMatcher*'s main prerequisite is written text.

B. Enhancing Bug Reports with User Feedback

Martens and Maalej [31] analyzed Twitter conversations between vendors' support accounts like @SpotifyCares and their users. Similarly to our statement, the authors highlight that users who provide feedback via social media are mostly non-technical users and rarely provide technical details. As support teams are interested in helping users, they initiate a conversation to ask for more context and details. They ask for context information like the affected hardware device, the app version, and its platform. Their objective is to better understand the issue to potentially forward that feedback to the development team and provide more helpful answers. Hassan et al. [14] show that developers also communicate with their users in the app stores to better understand their users.

Zimmermann et al. [53] show that the most important information in bug reports are steps to reproduce, stack traces, and test cases. The participants of their survey found that the version and operating system have lower importance than the previously mentioned information. However, the authors also argue that these details are helpful and might be needed to understand, reproduce, or triage bugs. Nevertheless, the authors did not focus on apps but developers and users of Apache, Eclipse, and Mozilla.

Developers could further use *DeepMatcher* to understand the popularity of bugs. They can achieve this in two steps. First, change *DeepMatcher* to take bug reports as an input to suggest problem reports (inverting the order as reported in the approach). Second, the parameter for the number of suggestions can either be increased or removed to enable suggesting all problem reports sorted by the similarity to the given bug report. This leads to an aggregated crowd-based severity level, a bug popularity score, or an indicator of how many users are affected by a certain bug report.

We further envision extracting context information and steps to reproduce from user feedback to enhance the issue tracker's bug report description. Having this information at hand can help developers narrow down the location of an issue and understand how many users are affected. Developers can use *DeepMatcher* to find problem reports related to bug reports by simply using a bug report summary as the input in our approach. Then, developers can skim through the suggested problem reports, select those that seem relevant, and then check whether they contain relevant context information. In case users did not provide useful information, developers can

take the IDs of the relevant problem reports and request more information from users in the Google Play Store. This process can partly be automated, e.g., using bots.

C. Extending DeepMatcher to Identify Duplicated, Recurring, or Similar Bug Reports

In Section IV, we found that *DeepMatcher* identified recurring bug reports. The Signal example in Table III shows a recurring bug report. Within the three bug report suggestions, *DeepMatcher* found three relevant matches. While the first bug report was filed in September 2015, the second in January 2016, the third in April 2016, a user reported the problem again in October 2017.

Consequently, developers might want to adapt *DeepMatcher* to either find recurring, similar, or duplicated bug reports even though it is not *DeepMatcher*'s primary goal. However, since the approach evaluates the matches based on context-sensitive text similarity, it could lead to promising results. Developers interested in these cases could, for example, increase *DeepMatcher*'s parameter *number of matching bug report suggestions* and use a bug report summary as the input for *DeepMatcher* to identify these cases. Future work could investigate and evaluate the use of *DeepMatcher* for such cases by utilizing our replication package.

VI. THREATS TO VALIDITY

We discuss threats to internal and external validity. Concerning the internal validity, we evaluated *DeepMatcher* by manually annotating 600 suggested bug reports for 200 problem reports. We performed two annotation tasks. One task to verify that the automatically classified app reviews are problem reports, and one to annotate whether *DeepMatcher*'s suggested matches are relevant for developers. As in every other manual labeling study, human coders are prone to errors. Additionally, their understanding of "a relevant match" may differ, which could lead to disagreements. To mitigate this risk, we designed both annotation tasks as peer-coding tasks. Two coders, each with several years of app development experience, independently annotated the bug report matches. For the verification of problem reports, we used a well-established coding guide by Maalej et al. [27], which Stanik et al. [44] also reused for the automatic problem report classification. To mitigate the threat to validity regarding the annotation of relevant matches, we performed test iterations on smaller samples of our collected dataset and discussed different interpretations and examples to create a shared understanding.

Further, we tried to collect a representative sample of meaningful app reviews. Thereby, we collected up to 10,000 app reviews for each app, ordered by helpfulness, covering more than two years. We did not aim for a comprehensive app review sample for a specific time frame but prioritized a meaningful app review sample from a larger time frame. Thereby, we could identify diverse insights within our qualitative analysis.

Another potential limitation is that we only considered 50 app reviews per app (200 in total), which we automatically

classified as problem reports. This classification might only find a specific problem report type, neglecting other informative problem reports. Other kinds of app reviews, including feature requests or praises, might also contain valuable information for developers, which *DeepMatcher* could match to bug reports. Therefore, our observations might differ for another sample of app reviews.

In the case *DeepMatcher* could not find any matching bug report among the three suggestions, we manually searched for relevant bugs in the issue trackers. We queried different term combinations and synonyms for certain features and components similar to how developers would proceed. However, not finding a relevant match in the issue tracker systems does not prove the non-existence of a relevant bug report in the issue tracker as we could have missed important terms in the query.

Concerning the external validity, our results are only valid for the four open source apps of our dataset. We considered different app categories, covering many tasks that users perform daily by including Firefox as an app for browsing the internet, Signal for messaging, Nextcloud for cloud storage, and VLC as a media player for music, videos, and streaming. However, these app categories include popular apps that we do not cover in our study, like Chrome or Safari. Further, the bug report suggestions could differ for closed source projects or apps of other mobile operating systems.

VII. RELATED WORK

A. User Feedback Analytics

Feedback-driven requirements engineering is an increasingly popular topic in research often focusing on *app reviews* [12], [13], [27], *tweets* [10], [49], product reviews such as *Amazon reviews* [23], [24], or a combination of reviews and product descriptions [19]. All of them have in common that a software product already exists and that users rate and write their experience with it after using it [37]. User feedback and involvement are important to both software engineers and requirements managers, as they often contain insights such as introduced bugs and feature requests [28], [44], [47]. The classification of user feedback [27] was a first step towards understanding user needs. Further studies [23], [24] looked at the classified feedback more closely by analyzing and understanding user rationale—the reasoning and justification of user decisions, opinions, and beliefs. Once a company decides to integrate, for example, an innovative feature request in the software product, it will be forwarded to the release planning phase [35], [36].

In our approach, we build on top of the existing body of research by, in particular, applying the machine learning approach of Stanik et al. [44] to identify problem reports in app reviews. We used that approach as an initial filter of the app reviews because Pagano and Maalej [37] showed that popular apps receive about 4,000 app reviews daily—which would be unfeasible for us to filter manually. Since the classification approach has an F1-score of 0.79 for identifying problem reports in English app reviews, we had to manually check the classified app reviews, as described in Section III-B.

B. Combining User Feedback and Bug Reports

El Mezouar et al. [7] present a semi-automatic approach that matches tweets with bug reports in issue tracking systems. They look at the bug reports of the two browsers Firefox and Chrome. They use natural language processing techniques to preprocess the text of both data sources and apply the Lucene search engine to suggest potentially matching bug reports. The approach crawls, preprocesses, filters, and normalizes tweets before they match them with issues. During the crawling process, the authors include tweets that either mention the browser with the @ symbol or hashtag. Then, they remove misspellings, abbreviations, and non-standard orthography. Afterward, the authors filter tweets with a list of bug related keywords like *lag* and *crash* while also considering negated bug terms with a part-of-speech analysis. In a final step, the approach removes symbols, punctuation, non-English terms, and stems the words using the porter stemmer [39]. For matching tweets with issues, the authors extract keywords from the tweets and use them as a search query in the Lucene search engine².

In contrast to El Mezouar et al., we consider app reviews from App Stores. While Tweets allow for lengthy conversations with stakeholders [11], [31] that may lead to in-depth insights into, e.g., the users' context like the app version and steps to reproduce, app stores enable developers to reply to app reviews, and users to update their review [14]. App reviews also contain metadata like the hardware device and the installed app version (that information is only available to the app developers). Further, stakeholders can ensure that users address the user wrote reviews for. When analyzing tweets and the software is available on multiple platforms like Windows, Mac, iOS, and Android, it is often difficult to understand which platform the user addressed without interacting with the user. Third, in app reviews, users can write longer texts than in tweets. Besides considering two platforms as our data source, we further applied more sophisticated technical solutions by applying state of the art NLP approaches that have a deeper understanding of the language than a search engine. We also build on top of previous research to extract problem reports from app reviews, leading to more relevant results than a simple keyword-based approach [27].

VIII. CONCLUSION

In this paper, we introduced *DeepMatcher*, an approach that extracts problem reports from app reviews submitted by users; and then identifies matching bug reports in an issue tracker used by the development team. Our approach primarily addresses the challenge of integrating user feedback into the bug fixing process. Developers may receive thousands of app reviews daily, which makes a manual analysis hard to unfeasible. Additionally, most user feedback is either praise like “I love this app.” or a dispraise like “I hate it!”. For the latter reason, we first filtered the problem reports from the reviews by reusing recent related work. After manually

²<https://lucene.apache.org>

validating the problem reports, we applied *DeepMatcher*, which takes a problem report and a bug report summary as the input. *DeepMatcher* then transforms the text into context-sensitive embeddings on which we applied cosine similarity to identify potential matching problem reports and bug reports. In total, from 200 problem reports, *DeepMatcher* was able to identify 167 relevant matches with bug reports. In 91 cases, *DeepMatcher* did not find any match. To understand whether indeed no match exists, we manually looked into corresponding issue trackers and found that in 44 cases, *DeepMatcher* missed a potential match while in 47 cases, no bug report existed. Our results show that our approach can help developers identify bugs earlier, enhance bug reports with user feedback, and eventually lead to more precise ways to detect duplicate or similar bugs.

ACKNOWLEDGMENT

This work was partly supported by the City of Hamburg (within the “Forum 4.0” project as part of the ahoi.digital funding line), by the European Union (within the Horizon 2020 EU project “OpenReq” Grant Nr. 732463), and by the German Federal Government (BMBF project “VentCore”).

REFERENCES

- [1] D. Ayata. Applying machine learning and natural language processing techniques to twitter sentiment classification for turkish and english.
- [2] D. Ayata, M. Saraclar, and A. Ozgur. BUSEM at SemEval-2017 task 4a sentiment analysis with word embedding and long short term memory RNN approaches. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 777–783. Association for Computational Linguistics.
- [3] D. Ayata, M. Saraclar, and A. Ozgur. Turkish tweet sentiment analysis with word embedding and machine learning. In *2017 25th Signal Processing and Communications Applications Conference (SIU)*, pages 1–4. IEEE.
- [4] L. V. G. Carreño and K. Winbladh. Analysis of user comments: an approach for software requirements evolution. In *2013 35th international conference on software engineering (ICSE)*, pages 582–591. IEEE, 2013.
- [5] F. Chollet. *Deep learning with Python*. Manning Publications.
- [6] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding.
- [7] M. El Mezouar, F. Zhang, and Y. Zou. Are tweets useful in the bug fixing process? an empirical study on firefox and chrome. *Empirical Software Engineering*, 23(3):1704–1742, 2018.
- [8] A. Finkelstein, M. Harman, Y. Jia, W. Martin, F. Sarro, and Y. Zhang. Investigating the relationship between price, rating, and popularity in the blackberry world app store. *Information and Software Technology*, 87, 2017.
- [9] J. Ghosh and A. Strehl. Similarity-based text clustering: A comparative study. In J. Kogan, C. Nicholas, and M. Teboulle, editors, *Grouping Multidimensional Data*, pages 73–97. Springer-Verlag.
- [10] E. Guzman, R. Alkadhi, and N. Seyff. A needle in a haystack: What do twitter users say about software? In *2016 IEEE 24th International Requirements Engineering Conference (RE)*, pages 96–105. IEEE, 2016.
- [11] E. Guzman, M. Ibrahim, and M. Glinz. A little bird told me: Mining tweets for requirements and software evolution. In *2017 IEEE 25th International Requirements Engineering Conference (RE)*, pages 11–20. IEEE, 2017.
- [12] E. Guzman and W. Maalej. How do users like this feature? a fine grained sentiment analysis of app reviews. In *2014 IEEE 22nd international requirements engineering conference (RE)*, pages 153–162. IEEE, 2014.
- [13] M. Harman, Y. Jia, and Y. Zhang. App store mining and analysis: Msr for app stores. In *Proceedings of the 9th IEEE Working Conference on Mining Software Repositories*, pages 108–111. IEEE Press, 2012.
- [14] S. Hassan, C. Tantithamthavorn, C.-P. Bezemer, and A. E. Hassan. Studying the dialogue between users and developers of free apps in the google play store. *Empirical Software Engineering*, 23(3):1275–1312, 2018.
- [15] M. Honnibal, I. Montani, S. Van Landeghem, and A. Boyd. spaCy: Industrial-strength natural language processing in python. Zenodo, 2020.
- [16] C. Iacob and R. Harrison. Retrieving and analyzing mobile apps feature requests from online reviews. In *2013 10th working conference on mining software repositories (MSR)*, pages 41–44. IEEE, 2013.
- [17] A. Islam and D. Inkpen. Semantic text similarity using corpus-based word similarity and string similarity. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 2(2):1–25, 2008.
- [18] V. Jijkoun, M. de Rijke, et al. Recognizing textual entailment using lexical similarity. In *Proceedings of the PASCAL Challenges Workshop on Recognising Textual Entailment*, pages 73–76. Citeseer, 2005.
- [19] T. Johann, C. Stanik, W. Maalej, et al. Safe: A simple approach for feature extraction from app descriptions and app reviews. In *2017 IEEE 25th International Requirements Engineering Conference (RE)*, pages 21–30. IEEE, 2017.
- [20] A. Joulin, E. Grave, P. Bojanowski, and T. Mikolov. Bag of tricks for efficient text classification. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 427–431, Valencia, Spain, Apr. 2017. Association for Computational Linguistics.
- [21] T. Kenter and M. De Rijke. Short text similarity with word embeddings. In *Proceedings of the 24th ACM international on conference on information and knowledge management*, pages 1411–1420, 2015.
- [22] A. J. Ko, B. A. Myers, and D. H. Chau. A linguistic analysis of how people describe software problems. In *Visual Languages and Human-Centric Computing (VL/HCC’06)*, pages 127–134. IEEE, 2006.
- [23] Z. Kurtanović and W. Maalej. Mining user rationale from software reviews. In *2017 IEEE 25th International Requirements Engineering Conference (RE)*, pages 61–70. IEEE, 2017.
- [24] Z. Kurtanović and W. Maalej. On user rationale in software engineering. *Requirements Engineering*, 23(3):357–379, 2018.
- [25] A. Lamkanfi, S. Demeyer, E. Giger, and B. Goethals. Predicting the severity of a reported bug. In *2010 7th IEEE Working Conference on Mining Software Repositories (MSR 2010)*, pages 1–10. IEEE, 2010.
- [26] H. Li, L. Zhang, L. Zhang, and J. Shen. A user satisfaction analysis approach for software evolution. In *IEEE International Conference on Progress in Informatics and Computing*, volume 2, 2010.
- [27] W. Maalej, Z. Kurtanović, H. Nabil, and C. Stanik. On the automatic classification of app reviews. *Requirements Engineering*, 21(3):311–331, 2016.
- [28] W. Maalej, M. Nayebi, T. Johann, and G. Ruhe. Toward data-driven requirements engineering. *IEEE Software*, 33(1):48–54, 2016.
- [29] C. Manning and H. Schütze. *Foundations of statistical natural language processing*. MIT press, 1999.
- [30] C. D. Manning, P. Raghavan, and H. Schütze. Chapter 8: Evaluation in information retrieval. *Introduction to information retrieval*, 10:1–18, 2009.
- [31] D. Martens and W. Maalej. Extracting and analyzing context information in user-support conversations on twitter. In *2019 IEEE 27th International Requirements Engineering Conference (RE)*, pages 131–141. IEEE, 2019.
- [32] D. Martens and W. Maalej. Release early, release often, and watch your users’ emotions: Lessons from emotional patterns. *IEEE Software*, 36(5):32–37, 2019.
- [33] S. McIlroy, N. Ali, and A. E. Hassan. Fresh apps: an empirical study of frequently-updated mobile apps in the google play store. *Empirical Software Engineering*, 21(3), 2016.
- [34] R. Mihalcea, C. Corley, C. Strapparava, et al. Corpus-based and knowledge-based measures of text semantic similarity. In *Aaai*, volume 6, pages 775–780, 2006.
- [35] M. Nayebi, L. Dicke, R. Ittype, C. Carlson, and G. Ruhe. Essmart way to manage user requests. *CoRR*, abs/1808.03796, 2018.
- [36] M. Nayebi and G. Ruhe. Asymmetric release planning-compromising satisfaction against dissatisfaction. *IEEE Transactions on Software Engineering*, 2018.
- [37] D. Pagano and W. Maalej. User feedback in the appstore: An empirical study. In *21st IEEE International Requirements Engineering Conference*. IEEE, 2013.
- [38] S. Panichella, A. Di Sorbo, E. Guzman, C. A. Visaggio, G. Canfora, and H. C. Gall. How can i improve my app? classifying user reviews

- for software maintenance and evolution. In *2015 IEEE International Conference on Software Maintenance and Evolution (ICSME)*, pages 281–290. IEEE, 2015.
- [39] M. F. Porter et al. An algorithm for suffix stripping. *Program*, 14(3):130–137, 1980.
- [40] X. Quan, G. Liu, Z. Lu, X. Ni, and L. Wenyin. Short text similarity based on probabilistic topics. *Knowledge and information systems*, 25(3):473–491, 2010.
- [41] V. Sanh, L. Debut, J. Chaumond, and T. Wolf. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter.
- [42] A. Simmons and L. Hoon. Agree to disagree: on labelling helpful app reviews. In *Proceedings of the 28th Australian Conference on Computer-Human Interaction*, pages 416–420, 2016.
- [43] P. Sitikhu, K. Pahi, P. Thapa, and S. Shakya. A comparison of semantic similarity methods for maximum human interpretability. In *2019 artificial intelligence for transforming business and society (AITB)*, volume 1, pages 1–4.
- [44] C. Stanik, M. Haering, and W. Maalej. Classifying multilingual user feedback using traditional machine learning and deep learning. In *27th IEEE International Requirements Engineering Conference Workshops*. IEEE, 2019.
- [45] C. Stanik, L. Montgomery, D. Martens, D. Fucci, and W. Maalej. A simple nlp-based approach to support onboarding and retention in open source communities. In *2018 IEEE International Conference on Software Maintenance and Evolution (ICSME)*, pages 172–182. IEEE, 2018.
- [46] Statista. Number of apps available in leading app stores as of 3rd quarter 2019. <https://www.statista.com/statistics/276623>. Accessed July 10, 2019.
- [47] L. Villarroel, G. Bavota, B. Russo, R. Oliveto, and M. Di Penta. Release planning of mobile apps based on user reviews. In *38th IEEE/ACM International Conference on Software Engineering*. IEEE, 2016.
- [48] X. Wang, L. Zhang, T. Xie, J. Anvik, and J. Sun. An approach to detecting duplicate bug reports using natural language and execution information. In *Proceedings of the 30th international conference on Software engineering*, pages 461–470, 2008.
- [49] G. Williams and A. Mahmoud. Mining twitter feeds for software user requirements. In *2017 IEEE 25th International Requirements Engineering Conference (RE)*, pages 1–10. IEEE, 2017.
- [50] G. Williams and A. Mahmoud. Modeling user concerns in the app store: A case study on the rise and fall of yik yak. In *26th IEEE International Requirements Engineering Conference*, 2018.
- [51] Y. Zhou, Y. Tong, R. Gu, and H. Gall. Combining text mining and data mining for bug report classification. *Journal of Software: Evolution and Process*, 28(3):150–176, 2016.
- [52] M. Zhu. Recall, precision and average precision. *Department of Statistics and Actuarial Science, University of Waterloo, Waterloo*, 2:30, 2004.
- [53] T. Zimmermann, R. Premraj, N. Bettenburg, S. Just, A. Schroter, and C. Weiss. What makes a good bug report? *IEEE Transactions on Software Engineering*, 36(5):618–643, 2010.