

Zero-shot Bilingual App Reviews Mining with Large Language Models

Jialiang Wei*, Anne-Lise Courbis*, Thomas Lambolais*,
Binbin Xu*, Pierre Louis Bernard** and Gérard Dray*

*: EuroMov Digital Health in Motion, Univ Montpellier, IMT Mines Ales, Ales, France

** : EuroMov Digital Health in Motion, Univ Montpellier, IMT Mines Ales, Montpellier, France

* : `firstname.lastname@mines-ales.fr` ** : `firstname.lastname@umontpellier.fr`

Abstract—App reviews from app stores are crucial for improving software requirements. A large number of valuable reviews are continually being posted, describing software problems and expected features. Effectively utilizing user reviews necessitates the extraction of relevant information, as well as their subsequent summarization. Due to the substantial volume of user reviews, manual analysis is arduous. Various approaches based on natural language processing (NLP) have been proposed for automatic user review mining. However, the majority of them requires a manually crafted dataset to train their models, which limits their usage in real-world scenarios.

In this work, we propose Mini-BAR, a tool that integrates large language models (LLMs) to perform zero-shot mining of user reviews in both English and French. Specifically, Mini-BAR is designed to (i) classify the user reviews, (ii) cluster similar reviews together, (iii) generate an abstractive summary for each cluster and (iv) rank the user review clusters. To evaluate the performance of Mini-BAR, we created a dataset containing 6,000 English and 6,000 French annotated user reviews and conducted extensive experiments. Preliminary results demonstrate the effectiveness and efficiency of Mini-BAR in requirement engineering by analyzing bilingual app reviews.

I. INTRODUCTION

App stores, such as Google Play and Apple App Store, allow users to express their feedback on downloaded apps. This feedback is in form of rating scores and text reviews. The latter contains praise and dispraise, user experience, problem reports, and feature requests [1]. App reviews are important for app success. As evidenced in prior literature, high-user rating scores have positive effects on apps’ sustainability [2]. Consequently, the design and development teams should take the users’ feedback into consideration during the evolution of their application.

Due to the large amount and the redundancy of app reviews, manual analysis is laborious. Various approaches based on natural language processing (NLP) have been proposed to reduce the efforts in analyzing user feedback, including the classification, clustering and summarization. Classification models are commonly employed in the first approach to categorize app reviews into predefined groups, such as feature requests and problem reports [3], [4], [5], [6], [7]. However, even after classification, the volume of reviews within each category remains substantial, making direct analysis impractical. To tackle this issue, some researchers have proposed grouping

reviews that pertain to the same topic [8], [9], [10], [11]. As the obtained clusters still contain a relatively large number of user reviews, the manual analysis of each cluster continues to be time-consuming. Certain techniques attempt to overcome this challenge by selecting the most representative phrases or sentences as summaries for groups of app reviews [9], [10], [12], [13], [14]. Nevertheless, this extractive summarization approach may not capture all the crucial information present within a given group. Moreover, existing approaches in user review analysis mainly focus on the English language, with few works on analyzing reviews in other languages [15], [16]. Furthermore, most existing approaches requires manually crafted dataset for training their classification models. The creation of dataset is costly and time consuming, which limits their usage in real-word scenarios. The objective of this article is therefore to address these gaps with large language models.

Pre-trained language models (PTMs) are deep neural networks previously trained on a vast corpus. Researchers have observed that large-sized PTMs display different behaviors from smaller PTMs and show surprising abilities (called emergent abilities) in solving a series of complex tasks. Thus, the research community coins the term “large language models (LLMs)” for these large-sized PTMs [17]. A remarkable application of LLMs is ChatGPT¹, it is fine-tuned from the GPT-3.5 [18] using Reinforcement Learning from Human Feedback (RLHF), which optimizes the model by interacting with human and learning from human preference. The Guanaco model is an open-source, finely-tuned LLM, derived through the application of QLoRa’s 4-bit tuning approach [19] on LLaMA base models [20]. QLoRA is an efficient fine-tuning approach that reduces memory usage. Guanaco is available in various parameter sizes, including 7B, 13B, 33B and 65B.

In this paper, we propose Mini-BAR, a bilingual approach based on LLMs to: (i) classify the user reviews into three categories: *feature request*, *problem report* and *irrelevant*; (ii) cluster similar reviews for *feature request* and *problem report*; (iii) generate a summary for each cluster of user reviews; and (iv) rank the user review clusters. Figure 1 depicts an overview of the workflow of Mini-BAR. We use the same pipeline to process the bilingual app reviews, eliminating the necessity of

¹<https://openai.com/blog/chatgpt/>

deploying separate models for each language. By combining these functionalities, Mini-BAR provides a comprehensive approach for analyzing bilingual app reviews, which can yield valuable insights for app developers and marketers. We validate the key steps of Mini-BAR by conducting an extensive set of experiments on 12000 annotated user reviews from three Health & Fitness apps. The results indicate that Mini-BAR has a satisfactory performance in both classification and clustering tasks, and produced high-quality summaries. We provide a replication package² containing the code, dataset, and experiment setups.

II. APPROACH

Mini-BAR provides support to developers for the analysis of mobile app user reviews through a four-step process. First, it applies a pre-trained classifier to categorize the user reviews (Section II-A). The second step clusters the user reviews based on their semantic similarity (Section II-B). The third step summarizes the user reviews belonging to the same cluster (Section II-C). The last step is to determine the importance of user review clusters and rank them accordingly (Section II-D). In the following subsections, we will detail each step of Mini-BAR.

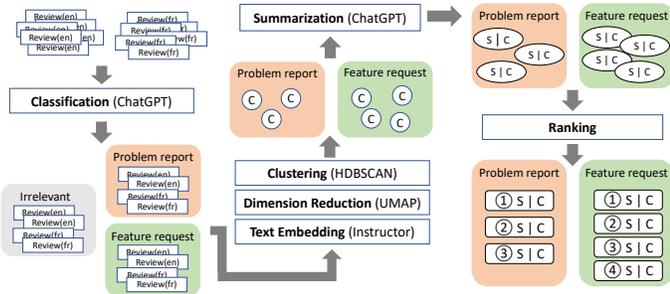


Fig. 1. Overview of Mini-BAR

A. Classification

The objective of this step is to automatically classify English and French user reviews into three categories: (F) *feature request*, (B) *problem report*, and (I) *irrelevant*. A user review may belong to either one of the three categories or both *feature request* and *problem report*. A user review is considered as *problem report* if it mentions the issues the users have experienced while using the app (e.g., “Can’t sync sleep data since last update”). *Feature requests* reflect users’ needs for new functions, new content, or improvements (e.g., “Please bring a feature to add some custom watch faces ...”). All the other user reviews are *irrelevant* (e.g., “Best app ever!”). Classifying the reviews can aid to redirect them to the appropriate software project members. For instance, *feature requests* can be delivered to requirements analysts, while *problem reports* can be directed to developers and testers [4].

²<https://github.com/JI-wei/mini-bar>

The classifier of Mini-BAR is based on ChatGPT, the model we use is *gpt-3.5-turbo*³. We use the following prompt to classify the app reviews.

```

Classify the following {lang} app review
into problem report, feature request or
irrelevant. Be concise.
...
{review}
...

```

Given a user review, its language is detected automatically with *Lingua*⁴, which is an accurate language detector. The detected language, which could be English, French among others, replaces the *{lang}* variable in the prompt. And the *{review}* in the prompt is replaced by the user review. The response of ChatGPT contains a single-phrase label name. We parse the response with regular expressions to automatically obtain the predicted labels.

B. Clustering

The objective of this step is to group English and French user reviews based on their semantic similarity, ensuring that reviews within a group are related to the same topic. Through clustering analysis, texts are divided into clusters such that those within a cluster exhibit semantic similarity. Currently, the RE community predominantly focuses on clustering English user reviews, leaving little attention to non-English user reviews. To address this gap, we propose a bilingual clustering approach that allows the creation of clusters comprising reviews from different languages that share common topics.

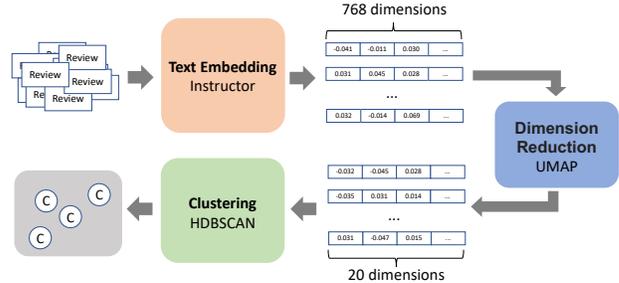


Fig. 2. Overview of clustering

In this step, we perform bilingual clustering analysis on user reviews that belong to the same category, namely *feature request* or *problem report*, which were identified in the previous step. It is worth noting that in the clustering process, the categories of *feature requests* and *problem reports* are processed separately, in order to obtain distinct clusters for each. Conversely, user reviews labeled as *irrelevant* are excluded from the clustering analysis.

The user reviews cannot be directly used as input for the clustering algorithm, as they are in a textual format. Therefore,

³<https://platform.openai.com/docs/models/gpt-3-5>

⁴<https://github.com/pemistahl/lingua-py>

it is necessary to convert them into embeddings — numerical vectors in a high-dimensional space — to enable effective processing. In this space, similar inputs in different languages are mapped close together. For example, the embeddings of “Problème de serveur récurrent” and “Connection issues to the main server” are in proximity to each other. We used *Instructor* [21] to embed English and French app reviews due to its high performance proven in Section III-B3. *Instructor* is an instruction-finetuned text embedding model that can generate text embeddings tailored to any task (e.g. clustering) by simply providing the task instruction, without any finetuning. In our case, we have utilized the instruction “Represent the app user review for clustering”. This model generates 768-dimensional embeddings for each app review.

The high dimension of embeddings causes high computation costs. Dimension reduction techniques can transform data from a high-dimensional space into a low-dimensional space and keeps meaningful information of the original data. As in Stanik et al. [9], we reduced the embeddings’ dimension with Uniform Manifold Approximation and Projection (UMAP). The implementation of Mini-BAR utilizes the UMAP Python package⁵, the UMAP parameters are as follows: output dimensionality of 20, number of neighbors set to 100, and minimum distance of 0.

The reduced embeddings of the *feature requests* and *problem reports* are then clustered by HDBSCAN [22], which has been proven efficient by Devine et al. [10] and Stanik et al. [9]. We use the HDBSCAN⁶ Python package for the implementation of Mini-BAR. The HDBSCAN parameters include a minimum cluster size of 5. The minimum cluster size is the smallest grouping size considered as a cluster. In our study, we chose a minimum cluster size of 5, as we were interested in identifying problems or features that were reported by at least 5 users.

C. Summarization

Given the potential magnitude of reviews within clusters, the process of summarization becomes imperative, enabling developers to efficiently grasp the cluster’s contents without the need to peruse every individual review. Large language models (LLMs), such as ChatGPT, has achieved a state-of-the-art performance in cross-lingual summarization [23]. In this step, we utilized ChatGPT (*gpt-3.5-turbo*) to generate abstractive English summaries for clusters containing bilingual user reviews. Our evaluation in Section III-C proves that ChatGPT outperforms extractive summarization method, and is able to generate high-quality summaries.

To generate the summaries, we utilized the following prompt within ChatGPT. The $\{reviews\ list\}$ is replaced by a list of user reviews separated by new line break. The response generated by ChatGPT represents the summary for that list of user reviews. Table I presents an example of a manual created cluster and its generated summary.

⁵<https://github.com/lmcinnes/umap>

⁶<https://github.com/scikit-learn-contrib/hdbscan>

```
Please summarize all following app
reviews into one English sentence:
```
{reviews list}
```
```

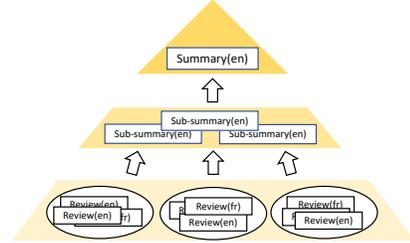


Fig. 3. Overview of summarization for large clusters

However, in case the user review clusters contain a large number of reviews and exceed the input length limitation of ChatGPT, it will issue a warning message indicating that “The message you submitted was too long, please reload the conversation and submit something shorter”. To address this issue, as illustrated in Figure 3, we adopted a hierarchical summarization approach consisting of the following steps: (i) dividing the reviews belonging to one cluster into multiple groups, each with a maximum of 4000 tokens⁷, (ii) generating a sub-summary for each group of user reviews, and (iii) obtaining the final summary by summarizing the sub-summaries. Since the length of all sub-summaries may also exceed the input limit of ChatGPT, we used a recursive procedure for steps (i) and (ii) by replacing the user reviews with sub-summaries.

TABLE I
EXAMPLE OF A USER REVIEW CLUSTER AND ITS GENERATED SUMMARY

User reviews:
<ul style="list-style-type: none"> - Dommage que la connexion 4g soit indispensable pour fonctionner. - Please for god sake make it to work offline also. - Is not work offline - It used to work offline. Now I have to log in just to see my old data. - Useless without internet.
Summary:
Users are disappointed that the app requires an internet connection to function and wish it could work offline like it used to.

D. Ranking

Given the clusters with summaries, the next step is to rank them by their importance. The goal of this step is to aid in the release planning of app developers. The importance of a cluster (*ClusterScore*) is determined based on the following characters:

- The number of reviews present within the cluster ($|reviews|$), problems or feature requests reported by

⁷<https://platform.openai.com/docs/guides/gpt/chat-completions-vs-completions>

more users should be given higher priority compared to those reported by fewer users.

- The average rating of the cluster (\overline{rating}). Clusters with lower average ratings should be given higher priority, as they may indicate users’ greater dissatisfaction with specific aspects of the app.
- The number of “thumbs up” inside the cluster ($|thumbsup|$). Users on Google Play have the option to click the “thumbs up” button on reviews that they find helpful. We posit that the number of “thumbs up” and the importance of a cluster are positively correlated as the review liked by more users should have a higher priority.

Given the weight of w_{rev} , w_{th} and w_{ra} , which are assigned default values of 1, 0.1, and 1, respectively, the calculation of *ClusterScore* is defined as follows:

$$ClusterScore = \frac{w_{rev} \cdot |reviews| + w_{th} \cdot |thumbsup|}{w_{ra} \cdot rating} \quad (1)$$

The clusters are ranked in decreasing order of *ClusterScore*.

III. EMPIRICAL EVALUATION

The objective of this study is to assess the performance of Mini-BAR with respect to three criteria: (i) its accuracy in classifying user reviews into one of three predefined categories, namely *feature request*, *problem report*, and *irrelevant*; (ii) its ability to cluster related user reviews that fall into the same topic; (iii) its ability to provide high-quality summaries of user reviews clusters. To achieve this goal, we evaluated Mini-BAR’s performance on a dataset of 6000 English and 6000 French reviews from three health-related mobile apps.

A. Evaluation of Classification

In this section, we aim to answer the following research question (RQ₁): *How accurate is Mini-BAR in classifying bilingual user reviews?*

1) *Dataset*: The training and evaluation of the classifier require a large number of labeled user reviews. We relabelled the 6000 French reviews of three applications (Garmin Connect, Huawei Health and Samsung Health) on Google Play from our previous work [16]. Besides, we collected 365,967 English user reviews from these three applications. For each application, 2000 English are randomly sampled for annotation. In this work, we have labeled 6000 English reviews and 6000 French reviews.

We used Prodigy⁸ from spaCy to annotate the user reviews. We created an annotation guide to clarify the definition of *feature request*, *problem report*, and *irrelevant*. Four authors of this paper annotated the sampled user reviews and they are finally reviewed by the first author of this paper. Table II shows the details of the annotated dataset. The sum of each category does not equal the total of reviews, as some reviews have been assigned to more than one label.

2) *Evaluation Metrics*: The performance of the classifiers is evaluated by *precision*, *recall*, and *F1* as presented in related work [15], [4].

⁸<https://prodi.gy/>

TABLE II
OVERVIEW OF THE DATASET FOR CLASSIFICATION

App	Language	Total	Feature request	problem report	Irrelevant
Garmin Connect	en	2000	223	579	1231
	fr	2000	217	772	1051
Huawei Health	en	2000	415	876	764
	fr	2000	387	842	817
Samsung Health	en	2000	528	500	990
	fr	2000	496	492	1047

3) *Experiments*: In this experiment, we compared the performance of ChatGPT and Guanaco-33B with ML models (Random Forest, Support Vector Machine), as well as various PTMs (BERT [24], CamemBERT[25], XLM-R [26]), on the classification of app reviews.

The ML models are trained using batch gradient descent, while PTMs employed mini-batch gradient descent, with a batch size of 12 and AdamW optimizer with a learning rate of $2e^{-5}$. They are trained on 3 epochs on a machine with a NVIDIA Tesla T4 GPU with 16 GB VRAM. The user reviews from the three apps of both languages were split using an 80:20 (training and test sets) ratio in a stratified manner, as illustrated in Figure 4. We trained the classifiers using a combination of *en_train* and *fr_train*. Subsequently, the classifiers were individually tested on the *en_test* and *fr_test*. We performed ten-fold cross-validation by randomly splitting the training and test sets ten times, and computed the average performance across these runs.

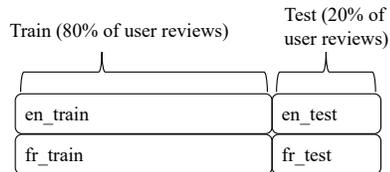


Fig. 4. Overview of dataset split for training and testing

We conduct classifications utilizing LLMs, specifically ChatGPT and Guanaco-33B, on all 12,000 user reviews. This is executed under zero-shot setting, implying that no prior training is involved. We have also assessed the performance of Guanaco-13B and Guanaco-65B. However, the responses generated by Guanaco-13B are disorganized, thereby hindering the extraction of predicted labels using regular expressions. The inference of Guanaco-65B is intolerably slow, even on advanced hardware such as the NVIDIA A100, making its usage impractical.

4) *Results*: The experiment results presented in Table III demonstrate that the ChatGPT exhibited good overall performance, which is comparable to that of ML models. Its comparatively lower performance in classifying *feature requests* can be attributed to the inherent complexity associated with such requests. In certain instances, users may express their desire for new features by criticizing existing ones or complaining about missing functionalities, rather than straightforwardly stating “I need...”. Among all the models, XLM-R archived

TABLE III
CLASSIFICATION ACCURACY ON USER REVIEWS OF THREE APPS

	Feature Request			Problem Report			Irrelevant			Average Weight			
	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1	
En	Random Forest	0.75	0.453	0.564	0.797	0.82	0.808	0.898	0.885	0.891	0.837	0.782	0.802
	SVM	0.86	0.438	0.58	0.86	0.806	0.832	0.931	0.893	0.912	0.895	0.778	0.823
	BERT	0.814	0.782	0.797	0.897	0.914	0.905	0.972	0.954	0.963	0.918	0.909	0.913
	CamemBERT	0.811	0.743	0.775	0.883	0.894	0.888	0.966	0.951	0.958	0.91	0.893	0.901
	XLM-R	0.823	0.811	0.816	0.902	0.917	0.909	0.979	0.958	0.968	0.925	0.917	0.92
	ChatGPT	0.768	0.5	0.606	0.762	0.972	0.854	0.983	0.911	0.945	0.871	0.853	0.852
	Guanaco-33B	0.361	0.62	0.456	0.662	0.951	0.781	0.983	0.817	0.893	0.763	0.823	0.774
Fr	Random Forest	0.8	0.528	0.635	0.798	0.834	0.816	0.902	0.869	0.885	0.848	0.796	0.817
	SVM	0.895	0.459	0.606	0.86	0.828	0.844	0.956	0.89	0.922	0.912	0.791	0.838
	BERT	0.766	0.725	0.744	0.871	0.866	0.869	0.947	0.931	0.939	0.888	0.872	0.88
	CamemBERT	0.852	0.823	0.837	0.922	0.925	0.923	0.977	0.96	0.968	0.936	0.924	0.929
	XLM-R	0.819	0.833	0.825	0.917	0.921	0.919	0.982	0.949	0.965	0.93	0.919	0.924
	ChatGPT	0.853	0.473	0.608	0.782	0.973	0.868	0.978	0.935	0.956	0.888	0.866	0.863
	Guanaco-33B	0.296	0.576	0.391	0.624	0.97	0.759	0.985	0.756	0.856	0.737	0.797	0.739

the best performance in bilingual classification. Although the performance of ChatGPT does not match up to PTMs, it is noteworthy that these LLMs have achieved such performance without the utilization of any reviews during their training phase. This suggests that ChatGPT can achieve satisfactory performance in user reviews of other application categories.

B. Evaluation of Clustering

In this section, our objectives are to address the RQ_2 : *How semantically meaningful are the clusters generated by MiniBAR?*

1) *Dataset*: In order to evaluate the performance of clustering algorithms, we created a dataset with 1200 user reviews. We randomly selected 100 problem reports and 100 feature requests from each of the three apps in each of the two languages present in the dataset created in Section III-A1. Then the authors employed manual clustering for each collection of 200 bilingual reviews, all of which pertained to an identical category. Reviews sharing the same topic were subsequently grouped into a single cluster. In instances where a user’s review mentioned multiple topics, the assignment of the cluster was determined by the initial topic that was reported. Two authors independently performed the clustering of the 1200 reviews, and their individual results were later merged through discussion. The resulting clusters were then considered as the ground truth for subsequent evaluation.

Table IV shows the number of manually created clusters and the number of clusters whose size is greater than or equal to 5 in each category and app. The feature requests encompass enhancements such as the modification of a Graphical User Interface (GUI), support for additional languages, increased activity options, integration with other applications, customization of permissions, and improvement of sleep tracking function. The problem reports predominantly center on a series of issues, notably the application’s unexpected crashes, errors encountered during the login process, difficulties in pairing with smartwatches, challenges with data synchronization, inconsistencies in the notification system, and complications related to Bluetooth connectivity, among others.

TABLE IV
OVERVIEW OF MANUALLY CREATED CLUSTERS

Bilingual	Garmin Connect	Huawei Health	Samsung Health
#clusters in feature request	89	74	69
#clusters($size \geq 5$) in feature request	7	9	11
#clusters in problem report	45	44	41
#clusters($size \geq 5$) in problem report	10	13	12

2) *Evaluation Metrics*: Following previous work [11], we use two commonly used indices, *Normalized Mutual Information (NMI)* [27] and *Adjusted Rand Index (ARI)* [28], to quantify the similarity between the automatic clustering and the ground truth. *NMI* ranges from 0 to 1, while *ARI* ranges from -1 to 1. A higher *NMI* or *ARI* indicates that the clustering method is more effective in producing clusters that align with the ground truth. Note that the *NMI* and *ARI* are computed for clusters with a size of 5 or greater, as the minimum cluster size of HDBSCAN is set to 5.

3) *Experiments and Results*: In this section, we evaluate the performance of different text representation methods (including traditional frequency-based methods, *bag of words (BOW)*, and *TF-IDF*, as well as PTM-based methods, *Universal Sentence Encoder* [29], *MiniLM* [30], *MPNet* [31], *E5* [32] and *Instructor* [21]) on the dataset created in Section III-B1. We performed three distinct experiments on English-only, French-only and bilingual user reviews. In each experiment, clustering is executed separately for the two categories (*feature requests* and *problem reports*) within each of the three applications.

4) *Results*: The average *NMI* and *ARI* of the three distinct experiments are shown in Table V. Results show that PTM-based methods outperformed traditional methods. Among all text representation methods, *Instructor* demonstrated the highest level of performance in clustering. The results from the English user reviews clustering were superior to those derived from the bilingual and French user reviews clustering. We attribute this variation to the relatively small French corpus employed for training the PTMs.

While the proposed approach appears to have some degree of validity, the results do not appear to be particularly encouraging at this stage. The primary reason for this is that

TABLE V
EVALUATION ON USER REVIEWS CLUSTERING

Embedding Methods	NMI			ARI		
	en	fr	bi	en	fr	bi
BOW	0.450	0.405	0.417	0.191	0.155	0.103
TF-IDF	0.452	0.449	0.460	0.253	0.216	0.149
USE	0.575	0.501	0.552	0.337	0.330	0.236
MiniLM	0.548	0.567	0.541	0.323	0.380	0.219
MPNet	0.616	0.575	0.593	0.400	0.346	0.278
E5	0.465	0.401	0.436	0.248	0.190	0.139
Instructor	0.713	0.587	0.603	0.597	0.357	0.308

our evaluation dataset is relatively small. Clustering algorithms require a sufficient amount of data to discover underlying patterns and structures. With a limited amount of data, it becomes difficult to identify meaningful groupings of text. Moreover, text clustering is a challenging task, particularly given the informal nature of the terminology employed in user reviews and the prevalence of spelling errors. Additionally, based on our empirical analysis, it appears that longer user reviews tend to result in less accurate clustering. This presents an opportunity for potential improvement by applying sentence-level clustering to user reviews.

C. Evaluation of Summarization

In this section, we aim to investigate the RQ_3 : *How effectively does ChatGPT perform in summarizing bilingual user reviews?*

1) *Dataset*: In Section III-B1, we carried out manual clustering for a total of 1200 user reviews. Among these, we utilized clusters with a size of 5 or greater to assess the performance of ChatGPT on summarization.

2) *Evaluation Metrics*: As outlined in Fabbri et al. [33], human evaluators rate the generated summaries based on four dimensions: *relevance* (the degree to which crucial information from the source has been included), *consistency* (how well the summary aligns with the factual details of the source), *fluency* (the quality of individual sentences), and *coherence* (the overall quality and coherence of all the sentences in the summary). Each dimension is scored on a Likert scale ranging from 1 to 5, with higher scores indicating superior performance.

3) *Experiments*: To evaluate the proficiency of ChatGPT in producing succinct English summaries of bilingual app reviews, we compare it with baseline approaches: Extractive summarization of Devine et al. [10], which selects the most representative sentence from a cluster of app reviews as the summary, the sentence is chosen by calculating the similarity with all other sentences of that cluster. Abstractive summarization with Guanaco models (13B, 33B and 65B version of Guanaco are used in our experiments) [19]. ChatGPT and Guanaco were instructed to synthesize clusters of user reviews into a single English sentence with the same prompt, as presented in Section II-C. Subsequently, human evaluators assessed the generated summaries. Given the pairs of user reviews and corresponding summaries, two authors of this paper were asked to evaluate the summaries on the Likert scale in the four dimensions that were previously mentioned.

TABLE VI
HUMAN EVALUATION ON GENERATED SUMMARIES

	Relevance	Consistency	Fluency	Coherence
Devine et al. [10]	4.23	4.83	4.62	4.71
Guanaco-13B	3.67	3.68	4.92	4.91
Guanaco-33B	4.65	4.58	4.91	4.88
Guanaco-65B	4.79	4.77	4.95	4.94
ChatGPT	4.81	4.84	4.95	4.94

4) *Results*: Table VI presents the average results of hand-made evaluations. Results show that abstractive approaches (ChatGPT, Guanaco-33B, and Guanaco-65B) generate very high quality sentences, and they are highly coherent when viewed in conjunction with one another. On the other hand, the extractive approach excels at extracting the most significant information from a cluster; however, it falls short in capturing all the essential details. During our manual evaluation, we found that Guanaco-13B tends to retrieve all available information available in the cluster without selectively focusing on crucial elements. In contrast, Guanaco-33B performs significantly better in this regard by effectively filtering out non-essential information. The results highlight the impressive performance of ChatGPT and Guanaco-65B in generating highly satisfactory summaries of user reviews.

IV. THREATS TO VALIDITY

This section aims to identify potential threats to the validity of our study.

1) *App reviews from one category*: All reviews studied in this paper are collected from three Health & Fitness apps (Garmin Connect, Huawei Health, and Samsung Health), mainly due to the context of health activity monitoring project. Instead of analyzing small number of reviews in many apps, we choose to annotate 12,000 reviews on three apps to create a larger dataset to evaluate clustering and summarization. However, these three apps may not be representative of apps in other categories. In the future, we will alleviate this threat by investigating user reviews of apps in various categories.

2) *Subjectivity in manual annotation*: The annotation of user reviews is a straightforward task, people without specific training can well classify or cluster the reviews. However, subjectivity can still arise during manual annotation, leading to variations in how different annotators interpret and label the reviews. To mitigate this threat, we (i) created a annotation guideline to detail the definition of each label following with examples, (ii) reviewed the final label through discussion and consensus.

3) *Issues of using ChatGPT*: ChatGPT has been selected as the classification and summarization component of Mini-BAR, owing to its superior capabilities. However, it is noteworthy to mention that certain countries have imposed prohibitions on the use of ChatGPT. Some users may refrain from utilizing ChatGPT owing to concerns pertaining to data privacy. And the cost of analyzing user reviews using ChatGPT cannot be overlooked, particularly in light of the large volume of user reviews. To mitigate these challenges, we have implemented alternative strategies. For classification, Mini-BAR users can

utilize our XLM-R checkpoint, which has been trained on 12,000 reviews, or they can opt for Guanaco-33B. For summarization tasks, users have the option to use other large language models, specifically Guanaco-33B or Guanaco-65B.

V. CONCLUSION

This paper introduces Mini-BAR, a mobile app review mining tool designed to assist app developers in extracting and summarizing user-reported issues and requests from a huge number of app reviews. This tool is based on LLMs and operates under zero-shot setting. Our empirical evaluation on the key steps of Mini-BAR resulted in numerous positive outcomes: (i) it accurately classified bilingual user reviews with an F1 score of 0.85; (ii) it created meaningful clusters of bilingual user reviews with a *NMI* greater than 0.6; (iii) it produced highly satisfactory summaries of bilingual user reviews.

In our future research, we intend to: (i) conduct a comparative analysis of various prompts utilized for classification and summarization, (ii) implement alternative large-scale language models for classification and summarization, (iii) execute classification and clustering at the sentence level as opposed to the review level, (iv) undertake evaluations using user reviews sourced from applications across a diverse range of categories.

REFERENCES

- [1] D. Pagano and W. Maalej, "User feedback in the appstore: An empirical study," in *2013 21st IEEE International Requirements Engineering Conference (RE)*, 2013, pp. 125–134.
- [2] G. Lee and T. S. Raghu, "Determinants of Mobile Apps' Success: Evidence from the App Store Market," *Journal of Management Information Systems*, vol. 31, no. 2, pp. 133–170, 2014.
- [3] N. Chen, J. Lin, S. C. Hoi *et al.*, "AR-miner: Mining informative reviews for developers from mobile app marketplace," *Proceedings - International Conference on Software Engineering*, no. 1, pp. 767–778, 2014.
- [4] W. Maalej, Z. Kurtanović, H. Nabil, and C. Stanik, "On the automatic classification of app reviews," *Requirements Engineering*, vol. 21, no. 3, pp. 311–331, 2016.
- [5] R. R. Mekala, A. Irfan, E. C. Groen *et al.*, "Classifying User Requirements from Online Feedback in Small Dataset Environments using Deep Learning," in *2021 IEEE 29th International Requirements Engineering Conference (RE)*, 2021, pp. 139–149.
- [6] P. R. Henao, J. Fischbach, D. Spies *et al.*, "Transfer Learning for Mining Feature Requests and Bug Reports from Tweets and App Store Reviews," in *2021 IEEE 29th International Requirements Engineering Conference Workshops (REW)*, 2021, pp. 80–86.
- [7] J. Zhang, Y. Chen, N. Niu, and C. Liu, "A Preliminary Evaluation of ChatGPT in Requirements Information Retrieval," no. 2022, pp. 1–16. [Online]. Available: <http://arxiv.org/abs/2304.12562> 2023.
- [8] S. Scalabrino, G. Bavota, B. Russo *et al.*, "Listening to the Crowd for the Release Planning of Mobile Apps," *IEEE Transactions on Software Engineering*, vol. 45, no. 1, pp. 68–86, 2019.
- [9] C. Stanik, T. Pietz, and W. Maalej, "Unsupervised Topic Discovery in User Comments," in *2021 IEEE 29th International Requirements Engineering Conference (RE)*, 2021, pp. 150–161.
- [10] P. Devine, J. Tizard, H. Wang *et al.*, "What's Inside a Cluster of Software User Feedback: A Study of Characterisation Methods," in *2022 IEEE 30th International Requirements Engineering Conference (RE)*, 2022, pp. 189–200.
- [11] Y. Wang, J. Wang, H. Zhang *et al.*, "Where is Your App Frustrating Users?" *Proceedings - International Conference on Software Engineering*, vol. 2022-May, pp. 2427–2439, 2022.
- [12] A. D. Sorbo, S. Panichella, C. V. Alexandru *et al.*, "What would users change in my App? Summarizing app reviews for recommending software changes," *Proceedings of the ACM SIGSOFT Symposium on the Foundations of Software Engineering*, vol. 13-18-Nove, pp. 499–510, 2016.
- [13] M. Alshangiti, W. Shi, E. Lima *et al.*, "Hierarchical Bayesian multi-kernel learning for integrated classification and summarization of app reviews," *ESEC/FSE 2022 - Proceedings of the 30th ACM Joint Meeting European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pp. 558–569, 2022.
- [14] C. Gao, Y. Li, S. Qi *et al.*, "Listening to Users' Voice: Automatic Summarization of Helpful App Reviews," *IEEE Transactions on Reliability*, pp. 1–13, 2022.
- [15] C. Stanik, M. Haering, and W. Maalej, "Classifying Multilingual User Feedback using Traditional Machine Learning and Deep Learning," in *2019 IEEE 27th International Requirements Engineering Conference Workshops (REW)*, 2019, pp. 220–226.
- [16] J. Wei, A.-L. Courbis, T. Lambolais *et al.*, "Towards a Data-Driven Requirements Engineering Approach: Automatic Analysis of User Reviews," in *7th National Conference on Practical Applications of Artificial Intelligence*, 2022. [Online]. Available: <http://arxiv.org/abs/2206.14669>
- [17] W. X. Zhao, K. Zhou, J. Li *et al.*, "A Survey of Large Language Models," pp. 1–85. [Online]. Available: <http://arxiv.org/abs/2303.18223> 2023.
- [18] T. Brown, B. Mann, N. Ryder *et al.*, "Language Models are Few-Shot Learners," in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell *et al.*, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 1877–1901.
- [19] T. Dettmers, A. Pagnoni, A. Holtzman, and L. Zettlemoyer, "QLoRA: Efficient Finetuning of Quantized LLMs." [Online]. Available: <http://arxiv.org/abs/2305.14314> 2023.
- [20] H. Touvron, T. Lavril, G. Izacard *et al.*, "Llama: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023.
- [21] H. Su, W. Shi, J. Kasai *et al.*, "One Embedder, Any Task: Instruction-Finetuned Text Embeddings." [Online]. Available: <http://arxiv.org/abs/2212.09741> 2022.
- [22] L. McInnes, J. Healy, and S. Astels, "hdbscan: Hierarchical density based clustering," *The Journal of Open Source Software*, vol. 2, no. 11, p. 205, 2017.
- [23] J. Wang, Y. Liang, F. Meng *et al.*, "Cross-Lingual Summarization via ChatGPT." [Online]. Available: <http://arxiv.org/abs/2302.14229> 2023.
- [24] J. Devlin, M. W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics*, vol. 1, oct 2019, pp. 4171–4186.
- [25] L. Martin, B. Muller, P. J. Ortiz Suárez *et al.*, "CamemBERT: a Tasty French Language Model," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 7203–7219.
- [26] S. Ruder, A. Søgaard, and I. Vulic, "Unsupervised cross-lingual representation learning," in *ACL 2019*, nov 2019, pp. 31–38.
- [27] N. X. Vinh, J. Epps, and J. Bailey, "Information theoretic measures for clusterings comparison: Is a correction for chance necessary?" *ACM International Conference Proceeding Series*, vol. 382, pp. 1073–1080, 2009.
- [28] L. Hubert and P. Arabie, "Comparing partitions," *Journal of Classification*, vol. 2, no. 1, pp. 193–218, 1985.
- [29] D. Cer, Y. Yang, S. yi Kong *et al.*, "Universal sentence encoder for English," *EMNLP 2018 - Conference on Empirical Methods in Natural Language Processing: System Demonstrations, Proceedings*, pp. 169–174, 2018.
- [30] W. Wang, F. Wei, L. Dong *et al.*, "MINILM: Deep self-attention distillation for task-agnostic compression of pre-trained transformers," *Advances in Neural Information Processing Systems*, 2020.
- [31] K. Song, X. Tan, T. Qin *et al.*, "MPNet: Masked and permuted pre-training for language understanding," *Advances in Neural Information Processing Systems*, vol. 2020-Decem, no. NeurIPS, pp. 1–14, 2020.
- [32] L. Wang, N. Yang, X. Huang *et al.*, "Text Embeddings by Weakly-Supervised Contrastive Pre-training," pp. 1–17. [Online]. Available: <http://arxiv.org/abs/2212.03533> 2022.
- [33] A. R. Fabbri, W. Kryściński, B. McCann *et al.*, "Summeval: Re-evaluating summarization evaluation," *Transactions of the Association for Computational Linguistics*, vol. 9, pp. 391–409, 2021.