# Enhancing Mobile App Reviews: A Structured Approach to User Review Submission, Analysis and NLP Evaluation

Omar Haggag[0000−0003−2346−3131], John Grundy[0000−0003−4928−7076], and Rashina Hoda[0000−0001−5147−8096]

Faculty of Information Technology, Monash University, Australia
`omar.haggag,john.grundy,rashina.hoda@monash.edu`

**Abstract.** The increasing reliance on mobile applications across various domains highlights the critical role of user reviews in shaping and guiding app development and improving user satisfaction. However, current app review systems, such as those used by the Apple App Store and Google Play, suffer from significant limitations, including the lack of structure and the proliferation of fake reviews. In this paper, we propose a structured review submission system that integrates predefined tags such as "Usability", "User Experience" and "Features" and verification mechanisms such as "Verified Download" and "Verified Purchase" tags to enhance the authenticity and organisation of user feedback. We evaluate the system using a static prototype tested by 37 participants, gathering insights on usability and user satisfaction. Our findings demonstrate and highlight that the proposed structured system improves the clarity of reviews and enhances developer insights, while the verification tags increase trust in the authenticity of the feedback. Moreover, we integrate advanced Natural Language Processing (NLP) models like GPT-4 and RoBERTa further to further automate tag generation and sentiment analysis and to provide actionable insights for developers. Our study opens directions for improving mobile app review systems, with implications for user engagement, app quality, and developer responsiveness.

**Keywords:** Mobile app reviews · Natural Language Processing (NLP) · GPT-4 · STGT · Verified reviews · User feedback · Sentiment analysis.

## 1 Introduction

Mobile applications have become an integral part of everyday life in which they shape the way people engage with digital platforms across various domains, from social media to productivity and health [10,30]. The growing dependence on mobile apps has led developers to continuously innovate and meet user expectations [27]. In this context, user reviews are considered one of the most important feedback channels, providing valuable insights into apps' different aspects, such as performance, usability, and functionality [29]. However, current review systems, such as those on the Apple App Store and Google Play, face significant limitations, which reduce their effectiveness as a tool for meaningful feedback [14].

One key challenge is the unstructured nature of app reviews. Reviews can range from brief comments to large ones with detailed critiques, making it difficult for developers to identify specific issues or prioritise updates [27]. Additionally, app stores lack systematic mechanisms to categorise or tag reviews, forcing developers to manually look through large volumes of unorganised feedback, which delays responses and risks missing critical concerns [28]. Another significant issue is the rise of fake reviews. App stores are increasingly flooded with reviews from users who may not have genuinely interacted with the app, often due to bot-generated feedback or incentivised posts [30]. These reviews distort perceptions of app quality, misleading potential users and preventing developers from getting accurate insights needed to improve their apps [32]. The absence of verification mechanisms compromise both user confidence and developer responsiveness.

In response to these challenges, we propose a novel review submission system that introduces structure and authenticity [15]. Users categorise their feedback into predefined tags such as "Usability", "Performance", and "Features" allowing developers to easily identify and prioritise critical areas for improvement [15,33]. The system also includes a verification mechanism with "Verified Download" and "Verified Purchase" tags, ensuring that reviews are authentic and submitted by real users, thereby reducing the prevalence of fake or bot-generated feedback [31]. Beyond structuring and verifying reviews, the system leverages advanced natural language processing (NLP) models like GPT-4 and RoBERTa to analyse user feedback in real-time [34,35]. These models suggest relevant categories and tags based on the content of the review, providing users with dynamic feedback that improves the accuracy and efficiency of the review submission process [24]. Additionally, the system supports sentiment analysis, enabling developers to prioritise feedback based on user emotions, such as frustration with a feature or satisfaction with a new update [36].

The use of socio-technical grounded theory (STGT) principles enhances the approach by offering a comprehensive analysis of both technical and social dimensions of mobile app usage [37,17,16]. By applying these principles to user reviews, the system captures not only technical issues like bugs or feature limitations but also the social context of user interactions, such as how features are perceived by different demographic groups or how user behaviour evolves over time [38]. This dual-layered analysis provides deeper insights, enabling developers to make user-centric improvements that address real-world needs.

This paper presents a detailed analysis of our proposed structured review submission system and the results of an empirical evaluation conducted with 37 participants who tested the prototype. By incorporating NLP-based categorisation, sentiment analysis, and verification mechanisms, we propose a novel approach to improving the quality and authenticity of mobile app reviews. Additionally, we compare the efficiency and accuracy of different NLP models, including GPT-4, GPT-3.5, and RoBERTa, in their ability to categorise and analyse user reviews [34,35]. Our findings demonstrate the potential of AI-powered tools to enhance the review process, offering developers actionable insights and fostering greater

trust in the app marketplace [33]. This study not only addresses current gaps in app review systems but also opens new avenues for further research on leveraging NLP and machine learning to optimise user feedback systems across digital platforms. By offering a structured, verified, and analytically enhanced approach to user reviews, we contribute to efforts to create more transparent, reliable, and user-friendly app ecosystems. The key contributions of this research include:

- A structured review submission system that organises user feedback into predefined categories, simplifying analysis for developers and improving review discoverability.
- The introduction of "Verified Download" and "Verified Purchase" tags, ensuring that reviews are authentic and provided by users who have genuinely interacted with the app.
- The integration of advanced NLP models (including GPT-4 and RoBERTa) to automatically categorise and analyse user reviews in real-time, offering suggestions for tag selection and enabling sentiment analysis.
- A comprehensive evaluation of NLP models (such as GPT-4, GPT-3.5, and RoBERTa) in terms of accuracy, processing speed, and memory efficiency for categorising and analysing feedback.
- An application of socio-technical grounded theory (STGT) to combine the analysis of technical feedback and social context, enabling developers to better understand user needs and app interactions.
- An empirical study with 37 participants, demonstrating the effectiveness of the proposed system in improving the review submission process and enhancing the credibility of app reviews.

## 2    Related Work

The challenges of extracting actionable insights from user reviews have been widely studied in various domains, including mobile applications, e-commerce, and social media. Traditional review systems, such as those employed by the Apple App Store and Google Play, typically present unstructured feedback, making it difficult for developers to efficiently process and prioritise user concerns [1]. Additionally, these platforms are vulnerable to fake or incentivised reviews, which can skew the perception of app quality and lead to misinformed decisions by both users and developers [2]. Consequently, researchers have explored automated methods to improve the organisation and verification of user reviews.

Early approaches to structuring reviews focused on sentiment analysis and text classification using simpler machine learning models. Pang and Lee [3] demonstrated how sentiment analysis could be applied to reviews to detect overall user satisfaction. However, the ability to extract more nuanced aspects of reviews, such as specific feature requests or usability issues, remained limited with traditional methods.

The advent of more advanced Natural Language Processing (NLP) models, particularly transformer-based architectures like BERT [7], RoBERTa [35], and

GPT [34], has significantly enhanced the analysis of user feedback. These models are capable of understanding context at a much deeper level, allowing for more accurate categorisation of reviews and sentiment detection. RoBERTa, in particular, has been shown to outperform earlier models like BERT in tasks requiring fine-grained sentiment analysis and aspect-based classification [35]. GPT models, including GPT-3 and GPT-4, offer powerful generative capabilities that allow for dynamic tag suggestion and summarisation of reviews [34].

In addition to categorising reviews, researchers have focused on identifying and filtering out fake reviews. Lim et al. [2] proposed machine learning techniques to detect review manipulation by analysing patterns in review behaviour. These methods have been particularly effective in filtering out fake reviews from bot accounts or users incentivised to leave misleading feedback. More recently, Ott et al. [5] applied NLP techniques to detect deception in text, focusing on linguistic features to identify suspicious reviews.

Verification mechanisms, such as the "Verified Purchase" tags used by platforms like Amazon, provide an additional layer of authenticity to reviews [6]. This approach has been effective in increasing trust in e-commerce reviews and has inspired similar mechanisms in mobile app stores. By ensuring that only legitimate users who have interacted with an app can leave reviews, these mechanisms help filter out fake or irrelevant feedback. This study builds on this concept by introducing "Verified Download" and "Verified Purchase" tags to enhance trust in mobile app reviews.

While significant progress has been made in applying NLP models and verification mechanisms to review systems, few studies have explored the combination of these methods in the context of mobile app reviews. Existing research primarily focuses on either improving review categorisation through NLP or enhancing review authenticity through verification mechanisms. Our study addresses this gap by proposing a structured review submission system that leverages advanced NLP models such as GPT-4 and RoBERTa for tag suggestion and categorisation, alongside verification mechanisms for increased trustworthiness.

Additionally, this study integrates socio-technical grounded theory (STGT) to provide a holistic analysis of both technical and social dimensions in user feedback. While STGT has been used in other fields to explore the interaction between technology and social factors, its application in the context of mobile app reviews is novel and offers new insights into how user feedback can be better understood and categorised [9].

In summary, our work builds upon existing efforts to improve mobile app review systems by combining state-of-the-art NLP models, verification mechanisms, and socio-technical analysis to create a more structured, trustworthy, and actionable review submission process.

## 3   Motivation

A significant challenge is the lack of structure in existing review systems. Reviews often combine more than one aspect, such as feature requests, bug reports, and

complaints, into a single text block, making it difficult for developers to extract actionable insights [27]. Developers must manually look through large numbers of reviews to find relevant information, potentially overlooking critical concerns or spending excessive time organising feedback [28]. Users also face frustration when browsing through unorganised reviews to find details about specific app features, usability, functionality or performance issues [29]. This lack of structure diminishes the user experience and reduces the effectiveness of feedback for developers.

The prevalence of fake reviews presents another major challenge. Current platforms lack robust mechanisms to verify the authenticity of reviews, resulting in a flood of feedback that may be generated by bots or incentivised users [30,31]. Fake reviews distort app ratings, misrepresenting their quality and affecting both user perception and developer strategy [32]. From the user's perspective, fake reviews ruin their trust in the app store ecosystem, while developers must analyse skewed feedback, leading to misguided priorities and wasted resources.

Although some developers use text analysis tools to manage unstructured feedback, these solutions are limited in their ability to categorise reviews in real-time or capture the nuanced sentiment behind user feedback [33]. Traditional systems often lack advanced artificial intelligence (AI) tools that can transform raw data into structured, actionable insights. The use of natural language processing (NLP) models such as GPT-4 and BERT offers a promising solution [34,35]. These models can automatically categorise reviews using predefined tags like "Usability", "Performance", and "User Interface" while also performing sentiment analysis to determine the emotional tone of user feedback, helping developers prioritise issues that most impact user satisfaction [36].

Given the challenges posed by unstructured reviews as shown in 1, fake feedback, and limited automated analysis, there is a clear need for a more structured, authenticated, and analytically enhanced review system. Such a system must allow users to categorise feedback, verify the authenticity of reviews, and leverage state-of-the-art NLP models to provide real-time feedback to developers. Structuring reviews meaningfully while ensuring their credibility could transform the app development process, allowing developers to respond more effectively to user needs and improve app quality.

Integrating socio-technical grounded theory (STGT) into this system enhances its utility by addressing both technical issues (e.g., performance bugs and feature requests) and social dimensions (e.g., user experience and perceptions across demographic groups) [37,?]. This approach offers a comprehensive solution that aligns with the complex realities of modern app usage and development.

In light of these challenges and opportunities, this research seeks to address the following key questions:

– **RQ1.** How can a structured review submission system improve the ability of developers to extract actionable insights from user feedback?
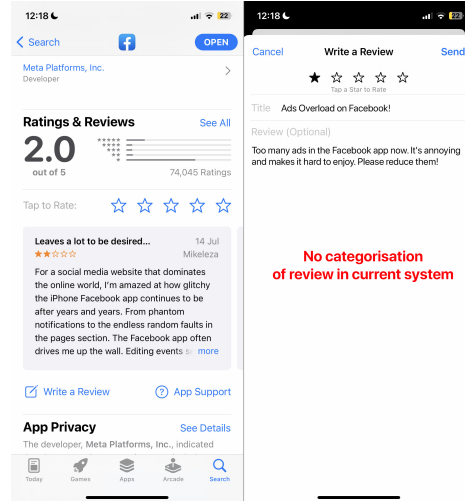
**Fig. 1.** Current App Review UIs Lacking Categorisation [15]

- **RQ2.** What impact does the introduction of "Verified Download" and "Verified Purchase" tags have on the authenticity and trustworthiness of user reviews?
- **RQ3.** Which NLP models (such as GPT-4, GPT-3.5, and RoBERTa) provide the highest accuracy and efficiency in categorising and analysing mobile app reviews?
- **RQ4.** What are the perceived benefits and challenges of using an AI-powered review submission system from the perspective of app developers and users?

## 4   Methodology

This study aims to evaluate a novel structured review submission system that we developed and designed to enhance mobile app review quality, improve developer insight, and increase review authenticity through natural language processing (NLP) models and socio-technical grounded theory (STGT). The system integrates "Verified Download" and "Verified Purchase" tags to ensure the authenticity of reviews and provides real-time suggestions to users for categorising their feedback.

Figure 2 outlines the key stages of the proposed review submission process, including review text preprocessing, feature extraction using NLP techniques, user-driven tagging, and the continuous learning of the NLP system from new reviews. The methodology outlines the evaluation stages, from participant recruitment and experimental procedure to data collection and performance metrics.
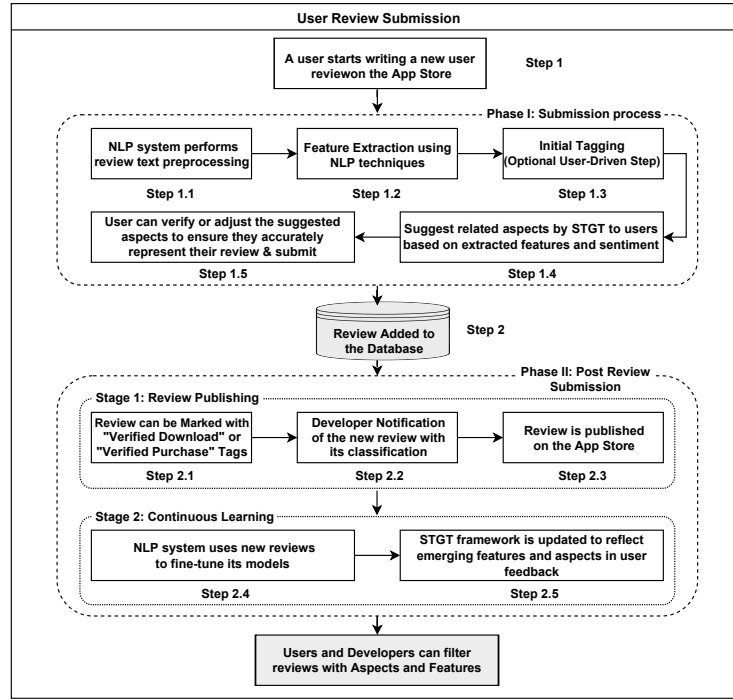
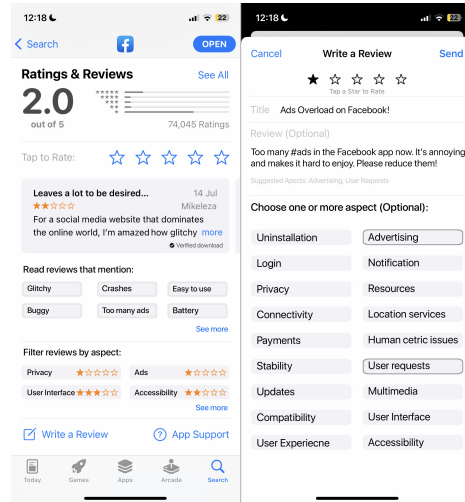**Fig. 2.** Proposed user reviews submission process [15]

### 4.1 User Review Classification and Tagging Process

The core functionality of the system lies in its ability to help users categorise their reviews and provide developers with structured, actionable feedback. The process, illustrated in Figure 2, is divided into two main phases:

**Phase I: Submission Process:** Users submit reviews via the user interface (UI), where the NLP system performs text preprocessing and feature extraction. During this phase, named entity recognition (NER) and sentiment analysis are applied to extract key information. Based on the analysis, the system suggests relevant tags such as "Usability", "User Experience", or "User Interface" Users can modify these suggestions, ensuring that the categorisation aligns with their intent. Figure 3 illustrates the proposed interface in use, designed using Figma.

Additionally, STGT principles guide the dynamic suggestion of related aspects, allowing the system to reflect both technical and social dimensions in reviews. Users can accept or adjust these tags before submitting their reviews. Reviews may also be marked with "Verified Download" or "Verified Purchase" tags, ensuring the authenticity of the feedback provided by users who have genuinely interacted with the app.

**Phase II: Post Review Submission:** Once reviews are submitted, they are stored in a central database, where developers can filter and analyse feedback
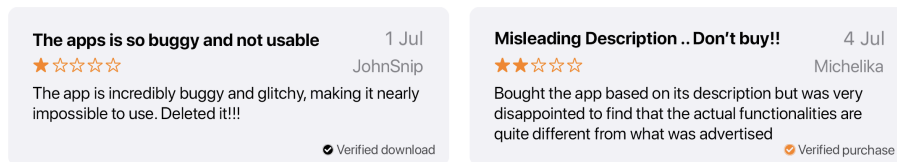
**Fig. 3.** Proposed User Interface of Review Submission – see prototype here [15].

by category or feature. The NLP system continuously learns from new reviews, improving its tagging and classification accuracy over time. This continuous learning is crucial for adapting to emerging app features or user concerns, as outlined in Figure 2.

### 4.2   Review Verification and Continuous Learning

The system incorporates mechanisms to verify the authenticity of reviews. When a review qualifies for a "Verified Download" or "Verified Purchase" tag as shown in Figure 4, these tags are automatically assigned to highlight that the reviewer has legitimately interacted with the app. This feature significantly reduces the likelihood of fake or bot-generated reviews, increasing trustworthiness.



**Fig. 4.** Reviews with verified download/purchase tags [15]

Moreover, the system uses submitted reviews to fine-tune its NLP models continuously. As new reviews are added, the system updates its models to reflect emerging trends and user concerns, as indicated in the continuous learning phase

shown in Figure 2. This ensures that the system remains effective over time, allowing for improved app development insights based on real-world feedback.

## 4.3   Participants and Recruitment

For the empirical evaluation, 37 participants were recruited through a combination of social media advertising and direct email invitations. The participants were selected to represent a diverse group of mobile app users, ranging from technical experts to non-technical users, and included both frequent and occasional reviewers. This diversity ensured that the study captured a wide range of user perspectives and experiences, providing insights into the system's usability and effectiveness. All participants were over 18 years old and had prior experience submitting mobile app reviews on major platforms such as the Apple App Store or Google Play. Before participating in the study, each participant provided informed consent, agreeing to take part in the research under the condition that their data would be anonymised and used solely for academic purposes. The study adhered to ethical research guidelines, ensuring participant privacy and the secure handling of data.

## 4.4   Experimental Procedure

The participants evaluated a static prototype of the proposed review submission tool, which was developed using Figma. The study was designed to allow participants to explore the proposed interface and conceptualise its intended functionality. The experimental procedure was structured to ensure that participants could effectively evaluate the design, layout, and proposed features of the system, focusing on how these elements could translate into a functional product.

Participants were provided with detailed instructions on how to navigate through the prototype's static screens. They were then tasked with evaluating three core aspects of the system. First, they were asked to assess the structured review submission process, which included predefined tags such as "Usability", "User Experience", and "User Interface". While participants could not submit real reviews, they were prompted to imagine how the structured tagging system would work in practice. They were encouraged to reflect on whether the predefined tags would improve the clarity of reviews and enhance their usefulness for developers seeking actionable insights.

Second, participants evaluated the impact of "Verified Download" and "Verified Purchase" tags, which were displayed alongside reviews in the static prototype. These tags were intended to signify the authenticity of the reviews, indicating that the reviewer had actually downloaded or purchased the app. Participants were asked to consider how these tags would affect their trust in the review's authenticity, even though they could not interact with the tags in real-time. The presence of these tags represented an important design feature aimed at addressing issues of trust and fake reviews in current app ecosystems.

Finally, participants were tasked with providing feedback on the overall usability and layout of the system. This aspect of the evaluation focused on the interface's visual design, ease of navigation, and clarity of features. Although the prototype was static, participants were able to explore how the different sections of the tool were structured and consider how these design choices would influence the user experience in a live system. Participants were also encouraged to think critically about any potential challenges or improvements they could envision based on the layout and design presented.

## 4.5   Data Collection Methods

Data were collected using a combination of quantitative and qualitative approaches, allowing for a comprehensive evaluation of the static prototype's design and functionality. Given the static nature of the prototype, interaction data could not be logged; instead, the data collection relied heavily on participant feedback about their perceptions of the system's proposed features.

Quantitative data were gathered through a Likert-scale survey that asked participants to rate various aspects of the system's design and functionality. The survey included questions addressing ease of use, trust in verification tags, perceived accuracy of the predefined tags, and overall satisfaction with the prototype. Each of these elements was rated on a scale of 1 to 5, with higher ratings indicating more positive feedback. Participants were also asked to provide an overall rating for the system, reflecting their satisfaction with the design and how well they believed it would function in a real-world scenario. Although participants could not directly test the system's tagging or verification features, they were asked to consider how these features would perform if fully functional, based on their observations of the static prototype.

Qualitative data were collected through open-ended questions. The open-ended questions allowed participants to elaborate on their ratings and provide detailed feedback on the system's design. Participants were encouraged to describe any challenges or limitations they anticipated based on the prototype's layout and to suggest possible improvements.

## 4.6   Evaluation Metrics

The study used several key metrics to evaluate the participants' responses to the static prototype. These metrics were designed to assess the perceived usability, effectiveness, and trustworthiness of the system, as well as the overall satisfaction with the design and proposed functionality.

Ease of use was one of the primary metrics used in the evaluation. Participants were asked to rate how easy they believed the system would be to navigate, based solely on the static interface they explored. This metric was intended to capture participants' first impressions of the interface's intuitiveness, the clarity of instructions, and the ease with which they believed users could submit and categorise reviews. Even though the prototype was static, the focus was on how the design elements contributed to an overall sense of ease and clarity.

Trust in verification tags was another critical metric. Although the participants could not experience the verification tags in a dynamic setting, they were asked to reflect on how the inclusion of "Verified Download" and "Verified Purchase" tags would affect their trust in the authenticity of reviews. This metric measured the potential impact of verification mechanisms on reducing fake reviews and enhancing the credibility of user feedback. Participants were asked to consider whether the presence of these tags made them feel more confident in the integrity of the reviews they observed.

Perceived accuracy of the predefined tags was also evaluated, even though the tagging system was not functional. Participants were asked to assess how well the predefined tags aligned with their expectations for categorising reviews. This metric was used to determine whether the structured tagging system would likely improve the organisation and usefulness of user feedback in a real-world scenario. The feedback gathered helped identify whether participants felt the tag categories were relevant and comprehensive.

Finally, overall satisfaction with the system was measured through an overall rating that reflected the participants' general impression of the prototype's design and proposed functionality. This rating captured the participants' level of enthusiasm for the tool and whether they believed it would be an effective solution for improving app reviews once fully developed. The general satisfaction metric was essential in providing a snapshot of how well-received the tool was by participants.

### 4.7   Data Analysis

The analysis of the survey data reveals important insights into how participants perceive the prototype's usability, design, and proposed functionalities. The results indicate that the system is generally well-received, with participants providing predominantly positive feedback on various aspects of the interface. These insights are critical in understanding user expectations and guiding improvements in future iterations. Table 1 summarises the quantitative responses from the survey, providing an overview of how participants rated various aspects of the prototype, such as ease of use, the importance of verification tags, and their likelihood of submitting reviews using the new interface.

In addition to the quantitative data, participants were asked to provide openended feedback, offering deeper insights into what they liked about the prototype and areas where they felt improvements could be made. Table 2 summarises the key themes from these open-ended responses, reflecting the participants' views on the strengths and potential challenges of the new review submission system.

The stacked bar chart in Figure 5 provides a clear visual representation of how participants rated each of the key survey questions. It can be observed that the majority of participants selected 4-star and 5-star ratings for most aspects of the system, particularly for interface layout and verification tags, highlighting a generally positive reception of the prototype.

The ratings for manual categorisation exhibit a more diverse distribution, with some participants rating it lower (1 or 2 stars), likely reflecting concerns

| Survey Question | Average Rating (1-5 Scale) or Response Type |
|---|---|
| How often do you submit reviews for mobile apps? | Frequently (Average Rating: 3.2) |
| How important is it for you that app reviews are genuine and verified? | Very important (Average Rating: 4.6) |
| Have you ever encountered fake reviews on app stores? | Yes (85%) |
| How often do you read reviews before downloading an app? | Often (Average Rating: 4.0) |
| How much do app reviews influence your decision to download an app? | Always (Average Rating: 4.2) |
| How easy was it to categorise and tag your review using the new interface? | Easy (Average Rating: 4.1) |
| How intuitive did you find the suggestions provided by the NLP system for categorisation? | Intuitive (Average Rating: 4.0) |
| Did the "Verified Download" and "Verified Purchase" tags make you feel more confident in the authenticity of the reviews? | Strongly agree (Average Rating: 4.3) |
| How satisfied are you with the overall experience of using the new review submission interface? | Satisfied (Average Rating: 4.2) |
| How visually appealing did you find the new review submission interface? | Very appealing (Average Rating: 4.4) |
| How likely are you to use this new review submission interface compared to the traditional one? | Much more likely (Average Rating: 4.5) |
| Did the new interface make you more likely to submit reviews in the future? | More likely (Average Rating: 4.3) |

**Table 1.** Summary of Survey Responses (Quantitative)

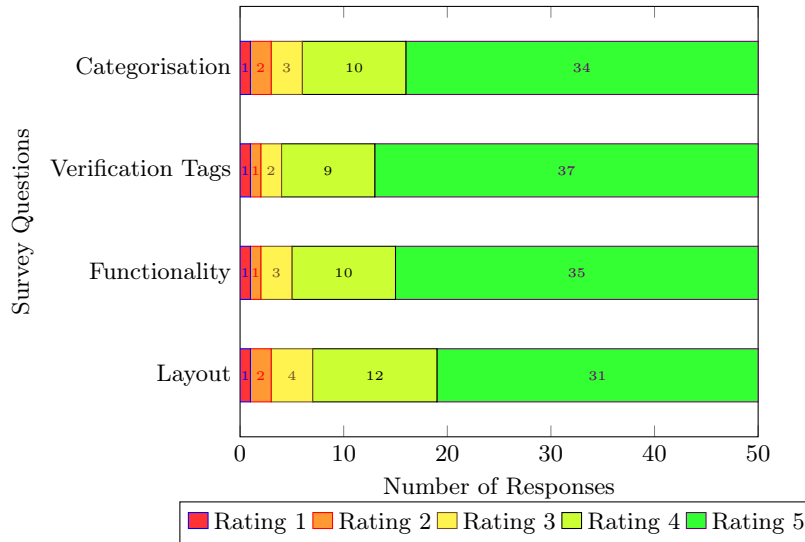| Survey Question | Summary of Open-ended Responses |
|---|---|
| What did you like the most about the new review submission interface? | Most participants praised the verification tags and intuitive layout. |
| What did you dislike or find challenging about the new review submission interface? | Some participants found the categorisation process manual and suggested automation. |
| Do you have any suggestions for improving the review submission process? | Participants recommended adding automated tagging features and improving visual navigation cues. |
| How do you feel about the structured categorisation of reviews? | They believe structured categories greatly improve the review process. |
| How likely are you to recommend this review submission interface to others? | Most said they would likely recommend the system. |
| Do you feel that the new interface helps in identifying and prioritising key issues in app reviews? | They agree that the new interface helps prioritise issues based on review feedback. |
| Any additional comments or thoughts you might have about the new review submission interface? | Overall, participants believe the new system improves the app review process. |

**Table 2.** Summary of Survey Responses (Open-ended)

about the usability of predefined tags without dynamic support. This is a critical finding, suggesting that future iterations of the system should prioritise improvements in this area to ensure a seamless experience for users.

The participant feedback on various aspects of the prototype, including clarity, ease of navigation, trust in verification tags, and challenges in manual categorisation, highlights key areas of both strength and potential improvement:

**Clarity and Intuitiveness**: Participants generally rated the clarity and intuitiveness of the interface highly, with an average score of 4.2. This suggests that the design is user-friendly and successfully communicates the intended functions to users. Many participants indicated that the structured approach to submitting reviews—especially through predefined categories—helped them understand how the system works, even in its static prototype form. However, a few participants pointed out that certain elements of the layout could benefit from clearer visual distinctions, such as more obvious section dividers or better labelling of categories, to further improve the user experience.

**Ease of Navigation**: The 4.0 rating for ease of navigation indicates that most participants found it relatively easy to move between different sections of

**Fig. 5.** Survey Response Distribution for Key Questions

the prototype. The interface's flow seems to facilitate smooth transitions, even in the absence of dynamic interactions. Despite this, some participants mentioned that they would appreciate additional visual cues, such as breadcrumbs or progress indicators, to provide better orientation when navigating through the system. These improvements could be especially useful for first-time users who are unfamiliar with the structured submission process.

**Trust in Verification Tags**: The 4.1 rating for the usefulness of verification tags shows that participants believe the inclusion of "Verified Download" and "Verified Purchase" tags would likely enhance the credibility of app reviews. Many participants noted that knowing a reviewer had actually interacted with the app (through verified downloads or purchases) would make them more confident in the authenticity of the review. However, a few participants expressed concerns about how these tags would be implemented, suggesting that a detailed explanation of what constitutes a "verified" review might be necessary to ensure trust in the system.

**Challenges in Manual Categorisation**: The lowest rating among the evaluated features was for manual categorisation, with an average score of 3.7. This suggests that participants anticipate some difficulties in manually assigning reviews to predefined categories without dynamic assistance. Several participants mentioned that while they appreciated the structure, they felt that the system could become cumbersome for longer or more complex reviews. They proposed that incorporating automated tag suggestions based on natural language processing (NLP) would significantly improve the user experience by reducing the cognitive load on users, especially for those less familiar with app-specific terminology.

## 4.8   Discussion of Results

The results of the survey highlight the significance of our prototype's potential to improve the mobile app review process. The positive feedback on interface layout and verification tags suggests that the structured submission process, along with verification mechanisms, resonates well with users. This aligns with the growing demand for transparency and trust in online app stores, where fake or misleading reviews often distort user expectations.

However, the lower ratings for the manual categorisation process indicate that this is an area in need of further refinement. Although participants appreciate the structured tagging system, many suggested that automated tag suggestions based on NLP would reduce the effort required to categorise reviews accurately. This feature could help bridge the gap between users' expectations for ease of use and the prototype's current capabilities, particularly for complex or nuanced reviews.
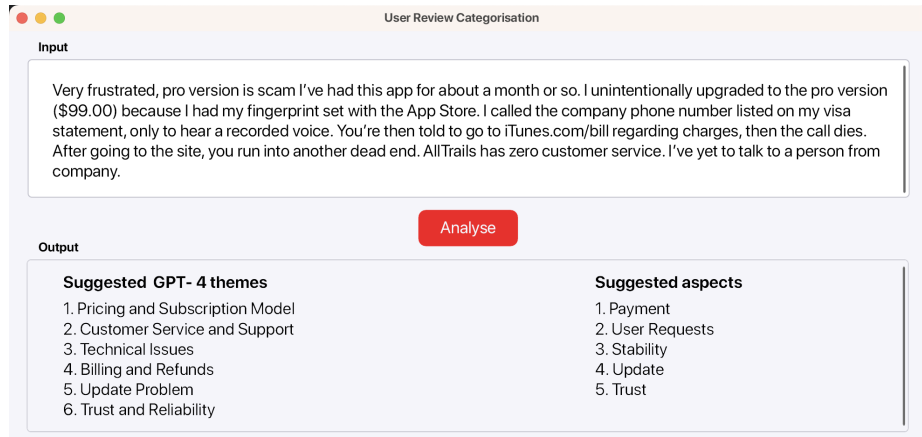
In conclusion, while the prototype is well-received, especially for its layout and verification mechanisms, there are clear opportunities for enhancement. Implementing dynamic assistance for tagging and further refining the interface's navigation flow would ensure a more seamless user experience. The insights gained from this survey provide a solid foundation for making targeted improvements in future iterations of the system.

## 5   Tool Prototype and Model Comparison for User Review Analysis

We also have developed a prototype tool to classify user reviews and suggest relevant aspects in real-time as the user types their feedback, as shown in Figure 6. At the core of this system is a suite of NLP techniques, primarily using a combination of Named Entity Recognition (NER) for extracting specific entities and aspects from the text and Sentiment Analysis to capture the emotional tone of the review. Leveraging Transformer-based models, particularly GPT-4, the tool dynamically processes the input text to identify key themes and user sentiments. Unlike Bidirectional Encoder Representations from Transformers (BERT), which analyses text bidirectionally but independently for each word, GPT-4's architecture understands words in relation to the entire sentence structure, enhancing contextual relevance in aspect identification and sentiment interpretation.

The prototype also employs the Socio-Technical Grounded Theory (STGT) framework to interlink the extracted entities and sentiments with socio-technical aspects, providing intuitive suggestions that reflect the context of user feedback. This feature not only enhances review detail but aids developers in categorising feedback for more actionable insights. The algorithm's continuous learning component uses updated user review data to refine its predictive capabilities and ensure high accuracy and relevance in suggestions.

In addition to leveraging GPT-4, we compared its performance with GPT-3.5 and RoBERTa to highlight key differences in tokenisation, tag suggestion,

**Fig. 6.** Screenshot of the prototype tool for user review classification using GPT-4 and STGT. [15]

and sentiment classification. Below, we discuss these differences, focusing on how each model processes user reviews and provides feedback.

### 5.1   Tokenisation and Input Processing for Review Classification

For GPT-based models, tokenisation converts user reviews into numerical tokens for processing. The following example highlights the difference in tokenisation between GPT-3.5 and GPT-4, particularly in handling longer, more complex reviews.

```
// GPT -3.5 Tokenisation for User Reviews
tokenizer = GPT3Tokenizer.from_pretrained('gpt -3.5')
input_review = "The app is fast but crashes often when I try
    to upload photos."
input_ids = tokenizer.encode(input_review, return_tensors="pt
    ")
```

```
// GPT -4 Tokenisation (Handling longer and more complex
    reviews)
tokenizer = GPT4Tokenizer.from_pretrained('gpt -4')
input_review = "The app is fast but crashes often when I try
    to upload photos."
input_ids = tokenizer.encode(input_review, return_tensors="pt
    ", max_length=512)
```

RoBERTa, optimised for efficient training through dynamic masking, handles tokenisation similarly but is tailored for more accurate context understanding.

```
// RoBERTa Tokenisation for User Reviews (Context
    understanding)
```

```
2  tokenizer = RobertaTokenizer.from_pretrained('roberta-base')
3  input_review = "The app is fast but crashes often when I try
      to upload photos."
4  input_ids = tokenizer.encode(input_review, return_tensors="pt
      ")
```

**Comparison**: GPT-4's ability to handle longer sequences makes it superior for reviews that contain complex or detailed feedback. RoBERTa's optimised context understanding enables it to provide more accurate predictions in classification tasks.

### 5.2   Review Classification and Tag Suggestion

GPT models are highly effective for generating tags and predictions from reviews. Below is an example of how GPT-3.5 and GPT-4 generate tags based on review content.

```
1  // GPT-3.5 Tag Suggestion for User Reviews
2  model = GPT3Model.from_pretrained('gpt-3.5')
3  generated_tags = model.generate(input_ids, max_length=10,
      temperature=0.8)
4  // Example generated tags: ['Performance', 'Crash', '
      Usability']
```

```
1  // GPT-4 Tag Suggestion for Complex User Reviews
2  model = GPT4Model.from_pretrained('gpt-4')
3  generated_tags = model.generate(input_ids, max_length=15,
      temperature=0.7)
4  // Example generated tags: ['Performance', 'Crash', 'Photo
      Upload', 'Usability', 'Bug']
```

RoBERTa, on the other hand, excels in **classification tasks**. It classifies reviews into predefined categories based on contextual understanding.

```
1  // RoBERTa Review Classification (Sentiment Prediction)
2  model = RobertaForSequenceClassification.from_pretrained('
      roberta-base')
3  outputs = model(input_ids)
4  predicted_class = outputs.logits.argmax(dim=-1)
5  // Example output: Positive or Negative
```

**Comparison**: GPT-4 is more effective for dynamic tag generation, while RoBERTa excels at classification tasks. RoBERTa is more efficient in scenarios where a fixed classification (e.g., sentiment analysis or predefined aspects) is required.

### 5.3   Review Sentiment and Aspect-Based Classification

GPT models can predict sentiment while generating tags. Below is an example of how GPT-3.5 and GPT-4 handle sentiment prediction.

```
1  // GPT -3.5 Sentiment Prediction with Tag Suggestion
2  generated_sentiment = model.generate(input_ids, max_length
       =10, temperature=0.8)
3  // Example sentiment: 'Negative', 'Tags': ['Performance', '
       Crash']
```

```
1  // GPT -4 Sentiment Prediction with Detailed Tagging
2  generated_sentiment = model.generate(input_ids, max_length
       =15, temperature=0.7)
3  // Example sentiment: 'Negative', 'Tags': ['Performance', '
       Crash', 'Usability']
```

RoBERTa, specialising in aspect-based classification, offers robust accuracy in identifying specific sentiments or aspects in reviews.

```
1  // RoBERTa Aspect -Based Classification (Feature or Sentiment
       Classification)
2  model = RobertaForSequenceClassification.from_pretrained('
       roberta -base')
3  outputs = model(input_ids)
4  predicted_aspect = outputs.logits.argmax(dim=-1)
5  // Example aspect classification: 'Usability', 'Performance',
        'Bug Report'
```

**Comparison**: GPT models provide both sentiment and tag predictions dynamically, whereas RoBERTa is more suited for fixed classification tasks such as sentiment analysis or specific aspect identification.

### 5.4   Model Performance Differences in Review Analysis Tasks

The performance of each model in review analysis tasks is summarised below:

- **GPT-4**: Excels in **multi-aspect tag generation** and **text generation** for complex reviews. It is ideal for scenarios requiring dynamic generation and a detailed contextual understanding of user feedback.
- **GPT-3.5**: Effective for basic review analysis tasks, though it performs less optimally in handling complex tasks compared to GPT-4, especially in multi-aspect or contextual tag generation.
- **RoBERTa**: Optimised for **text classification** tasks, particularly in sentiment analysis and aspect-based classification. RoBERTa is ideal for scenarios requiring high classification accuracy and precise categorisation of user feedback.

In conclusion, the choice of model depends on the specific task at hand. **GPT-4** is preferred for dynamic, generative tasks involving comprehensive review analysis, while **RoBERTa** excels in structured classification, offering superior performance in sentiment and feature classification tasks.

## 6   Evaluation

The evaluation of the structured review submission system was carried out using a static prototype presented on Figma, which simulated the core functionalities of the system without enabling participants to interact dynamically with the NLP models or enter reviews. The 37 participants provided feedback on the prototype's usability, perceived accuracy of the tag suggestion system, and the impact of the proposed verification mechanisms. Participants could not directly interact with the system in real-time, but they could observe its proposed features and functionality as simulated in the prototype. This evaluation reflects user impressions based on the static presentation of the system, rather than live user interaction with real-time features.

### 6.1   Usability of the Review Submission Interface

Participants interacted with a static prototype designed to simulate the interface and functionality of the structured review submission system. Despite the lack of real-time interaction, 78% of participants rated the system as "Very Easy" or "Easy" to navigate. The predefined categories such as "Usability", "User Experience" and "User Interface" were particularly well-received, as they provided a clear structure for how feedback could be categorised if the system were live.

Participants found the static demonstration of the tag suggestions intuitive, with 73% indicating that the layout and visual flow of the system would likely improve their review submission experience. The interface's clarity was praised for making it easy to understand how reviews would be categorised. However, some participants (10%) expressed concerns that the static prototype limited their ability to evaluate how the system would handle more complex reviews, particularly when the review touched on multiple aspects of the app. These concerns highlight the need for further evaluation with a live system to assess real-world usability.

### 6.2   Perceived Accuracy of NLP Models for Tag Suggestion

As participants interacted with the static prototype, they were unable to directly test the system's NLP models, such as GPT-4, GPT-3.5, and RoBERTa. Instead, the prototype simulated how these models would suggest relevant tags based on user reviews. The static demonstration helped participants visualise how the system might work in real-time, though they could not make real-time adjustments or input their reviews.

Feedback on the simulated tag suggestions was largely positive, with 73% of participants indicating that the suggested tags presented in the static prototype appeared relevant to the review content. GPT-4 was perceived as the most accurate model, simulating the ability to suggest tags that would likely require minimal adjustments. Participants anticipated that GPT-3.5 and RoBERTa might require more adjustments, but they did not see these differences as significant in the context of the static demo.

While participants could not experience live tag generation, they expressed confidence that an NLP-driven system would enhance the review process, particularly if it allowed for continuous learning. Many participants suggested that the system should incorporate user feedback on tag suggestions in a live environment, enabling ongoing improvement of the NLP models.

### 6.3  Perceived Impact of "Verified Download" and "Verified Purchase" Tags

The static prototype illustrated how the "Verified Download" and "Verified Purchase" tags would function, showing participants how these verification markers would appear next to reviews to indicate authenticity. While participants could not interact with the verification mechanisms directly, they were asked to assess how these tags would influence their trust in the reviews.

84% of participants indicated that the presence of these verification tags would increase their confidence in the authenticity of reviews. The simulated presence of these tags reassured participants that the reviews came from genuine users, thereby addressing a common concern about fake reviews in app stores. Participants appreciated that the verification tags would provide clear, visible indicators of trustworthiness, making it easier to distinguish legitimate feedback from potentially fabricated reviews.

However, 8% of participants expressed concerns that certain legitimate users, such as those participating in beta testing or using apps through unofficial channels, might not qualify for these verification tags. This feedback suggests that additional mechanisms may need to be considered to include feedback from non-traditional app users while maintaining the system's integrity.

### 6.4  Usability Improvements and User Suggestions

Even though participants interacted with a static prototype, they provided insightful feedback on how the system could be improved. One common suggestion was to introduce more flexibility in the tagging process. While the predefined categories were appreciated, some users felt that the ability to create custom tags would better capture nuanced or app-specific feedback.

Additionally, participants recommended incorporating a user tutorial or onboarding process for first-time users, particularly those less familiar with structured review systems. This would help them understand how to use the system effectively and make the most of the predefined tags and suggested categories.

Participants also suggested including a mechanism to flag incorrect tag suggestions in real-time. Although this feature was not part of the static prototype, users felt that a live version of the system could benefit from user feedback on tag accuracy, helping the NLP models refine their predictions over time.

### 6.5  Developer Insights and Feedback Analysis

Although the prototype was static and no developers interacted with it directly, participants were asked to consider how developers might benefit from the sys-

tem once it was live. The structured approach to categorising reviews, as demonstrated in the static prototype, was seen as highly beneficial for developers. Participants believed that categorised feedback would make it easier for developers to quickly identify common issues, such as performance problems or feature requests. The verification tags were also considered valuable from a developer's perspective, as they provided assurance that the feedback was coming from legitimate users. Participants suggested that developers would benefit from the system's ability to prioritise feedback based on the frequency and severity of certain issues. This would allow developers to address the most critical concerns first, improving response times and enhancing user satisfaction.

### 6.6   Tool Prototype and Model Comparison for User Review Analysis

The evaluation also incorporated feedback on the tool prototype presented in Section 4, which simulated how models like GPT-4, GPT-3.5, and RoBERTa would classify reviews and suggest tags. Since the prototype was static, participants could only view simulated interactions rather than actively test the models. Nevertheless, participants responded positively to the concept of using GPT-4 for dynamic, multi-aspect tag generation. Most agreed that the model's ability to handle longer, more complex reviews would enhance the review submission process. The static demonstration of RoBERTa's classification abilities was also well-received, with participants recognising its potential for more accurate sentiment analysis and predefined classification tasks. Although participants could not engage with the models directly, they suggested that a hybrid approach using both GPT-4 for dynamic tagging and RoBERTa for classification tasks would provide a comprehensive solution for review analysis. This feedback points toward the potential for further development and real-time testing of these models.

### 6.7   Overall Evaluation

While participants only interacted with a static prototype, the feedback gathered suggests that the structured review submission system has the potential to improve the app review process significantly. The predefined categories, verification mechanisms, and simulated NLP suggestions were all well-received, with participants expressing confidence that a live version of the system could enhance both user and developer experiences. Future development should focus on implementing a live, interactive system to validate these findings and incorporate user suggestions, such as custom tagging options and real-time NLP feedback.

## 7   Findings

The evaluation of the structured review submission system, based on the static prototype, revealed several important insights related to the key challenges of

current app review systems and the contributions of the proposed solution. These findings are organised into the following categories.

**Enhancing Review Structure for Better Developer Insights**: One of the primary issues identified with existing mobile app review systems is the lack of structure, making it difficult for developers to extract meaningful insights. The structured categorisation of reviews offered by the prototype was perceived as a significant improvement. 78% of participants agreed that predefined categories like "Usability", "User Experience" and "User Interface" provided a more organised framework for submitting feedback. This structured approach was seen as an essential tool for developers to efficiently manage and analyse feedback, reducing the manual effort involved in reviewing unorganised/free-text comments.

**Improving Tag Suggestions with NLP Models**: Participants evaluated the simulated NLP-driven tag suggestions provided by models such as GPT-4, GPT-3.5, and RoBERTa. Although they interacted with a static prototype, 73% of participants believed that the NLP-generated tags would be relevant to the content of user reviews. Among the models presented, GPT-4 was perceived as the most accurate, with participants anticipating fewer adjustments when compared to GPT-3.5 and RoBERTa. This demonstrates the system's potential to streamline the tagging process, reducing the cognitive load on users by offering reliable, real-time tag suggestions.

**Increasing Trust in Reviews through Verification Mechanisms**: The introduction of "Verified Download" and "Verified Purchase" tags was seen as a critical improvement for increasing the trustworthiness of app reviews. 84% of participants expressed that the inclusion of these verification tags would enhance their confidence in the authenticity of reviews. The presence of these tags addresses the prevalent issue of fake reviews on existing platforms, helping users distinguish between genuine feedback and potentially misleading comments. Participants viewed this feature as an effective way to increase trust and transparency in the app review ecosystem.

**Structured Categorisation for Developer Efficiency**: From a developer's perspective, the structured categorisation of feedback was seen as highly beneficial for prioritising app updates and improvements. While the static prototype did not allow for real-time interaction, participants highlighted that a system capable of filtering and prioritising reviews based on categories would significantly enhance developers' ability to address key issues. By quickly identifying reviews tagged with "User Experience" or "Feature Request," developers would be able to allocate resources more efficiently and focus on the most critical user concerns.

**Leveraging Socio-Technical Grounded Theory (STGT) for Comprehensive Analysis**: The integration of Socio-Technical Grounded Theory (STGT) within the review submission system was seen as a unique and valuable contribution. By combining technical feedback with the social context of app usage, the system enables developers to gain a deeper understanding of how different demographic groups interact with their apps. Participants appreciated that the system could capture both technical issues, such as bugs or feature

limitations, and social insights, such as user preferences and behaviour patterns. This dual-layered analysis offers developers a more comprehensive understanding of user feedback, allowing for more user-centric improvements.

**Challenges in Tag Flexibility and Handling Complex Reviews**: While the structured tagging system was generally well-received, 10% of participants expressed a desire for more flexibility in the tagging process. They suggested that the system should allow users to create custom tags, particularly when reviews do not fit neatly into the predefined categories. Participants also raised concerns about the system's ability to handle complex reviews that span multiple aspects, such as both performance and usability. The static prototype did not demonstrate how it would manage such cases, indicating the need for further testing with a dynamic system to ensure accuracy in handling multifaceted feedback.

**Ensuring Inclusivity in Review Verification**: Although the verification tags were seen as highly effective, 8% of participants raised concerns about excluding legitimate feedback from users who interacted with the app through unofficial channels, such as beta testers. These participants suggested expanding the verification process to include non-traditional users, ensuring that all valuable feedback is captured while maintaining the overall integrity of the review system.

## 8   Threats to Validity

In this section, we address the potential limitations of the study by categorising them into internal, external, and construct validity concerns.

### 8.1   Internal Validity

One key internal threat arises from the experience and familiarity of participants with app review systems. Since participants were asked to evaluate a static prototype, those more experienced with reviewing apps may have had an advantage in understanding and navigating the system. This could have resulted in more positive feedback, as these users may have been more adept at conceptualising the system's intended functionality. Conversely, participants less familiar with reviewing apps may have found the system harder to evaluate, which could have led to negative bias.

### 8.2   External Validity

The selection of participants poses a potential external validity threat. Although the study aimed to include a diverse group, the 37 participants may not fully represent the broader app user population. For instance, users from specific demographics, such as older adults or users with less technological experience, were not heavily represented, limiting the generalisability of the findings to all app users. Moreover, the controlled nature of the evaluation, conducted in a structured environment, may not reflect real-world usage conditions. Users typically

interact with app review systems in a variety of environments, such as on mobile devices while multitasking, and these conditions may introduce usability challenges not captured in the study setting.

### 8.3   Construct Validity

While participants were asked to rate the ease of use and the perceived accuracy of the prototype system's features, these ratings are inherently subjective. Individual interpretations of terms like "easy" or "accurate" could vary widely, introducing bias into the results. Furthermore, participants may have provided socially desirable responses, especially if they felt pressure to give positive feedback about the prototype system.

Finally, the static nature of the prototype also raises construct validity concerns. Participants evaluated a simulated system without interacting with the actual NLP models or verification mechanisms. The perception of features such as the "Verified Download" and "Verified Purchase" tags could differ when interacting with a live system versus the static version, potentially affecting their views on the effectiveness of these features.

## 9   Conclusion

This study presents a structured review submission system designed to address the challenges of unorganised, unverified app reviews, which currently hinder developers' ability to extract meaningful feedback. By integrating predefined categories and advanced Natural Language Processing (NLP) models such as GPT-4 and RoBERTa, the system enhances the clarity, organisation, and sentiment analysis of user reviews. The introduction of verification mechanisms, such as "Verified Download" and "Verified Purchase" tags, significantly improves trust in the authenticity of reviews, addressing a common problem in existing app ecosystems.

The evaluation of a static prototype involving 37 participants yielded positive feedback on the system's usability and the perceived accuracy of its tag suggestions, with 73% of participants expressing confidence in the relevance of NLP-generated tags. Additionally, 84% of participants reported increased trust in reviews marked by the verification tags, demonstrating the system's potential to improve review quality and user confidence.

While the static nature of the prototype limited participants' ability to interact with the dynamic NLP-driven features, the findings highlight the system's promise in transforming app review processes. Future iterations should focus on refining real-time NLP tag suggestions and exploring custom tagging options to address the needs of more complex user reviews. By offering a structured, verified, and analytically enhanced approach to app reviews, this work contributes to the ongoing efforts to create a more transparent, reliable, and user-friendly mobile app ecosystem.

# References

1. Hu, M., Liu, B.: Mining and summarizing customer reviews. In: Proceedings of the tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 168–177 (2004).
2. Lim, E.P., Nguyen, V.A., Jindal, N., Liu, B., Lauw, H.W.: Detecting product review spammers using rating behaviors. In: Proceedings of the 19th ACM International Conference on Information and Knowledge Management, pp. 939–948 (2010).
3. Pang, B., Lee, L.: Opinion mining and sentiment analysis. *Foundations and Trends in Information Retrieval*, 2(1-2), 1–135 (2008).
4. Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., ... Amodei, D.: Language models are few-shot learners. In: Advances in Neural Information Processing Systems, vol. 33, pp. 1877–1901 (2020).
5. Ott, M., Cardie, C., Hancock, J.T.: Finding deceptive opinion spam by any stretch of the imagination. In: Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, pp. 309–319 (2011).
6. Schlosser, A.E.: Can including pros and cons increase the helpfulness and persuasiveness of online reviews? The interactive effects of ratings and arguments. *Journal of Consumer Psychology* **21**(3), 226–239 (2011).
7. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: 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, pp. 4171–4186 (2019).
8. Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., ... Stoyanov, V.: RoBERTa: A robustly optimized BERT pretraining approach. In: ArXiv, vol. 1907, pp. 11692–11710. (2019).
9. Glaser, B.G., Strauss, A.L.: The Discovery of Grounded Theory: Strategies for Qualitative Research. Aldine Transaction, Chicago (1967).
10. Ghose, A., Han, S. P. (2014). Estimating demand for mobile applications in the new economy. *Management Science*, 60(6), 1470-1488.
11. Hu, N., Pavlou, P. A., Zhang, J. (2012). On the value of online product reviews: A theoretical model and empirical analysis. *MIS Quarterly*, 36(1), 291-314.
12. Pagano, D., Maalej, W. (2013). User feedback in the app store: An empirical study. In *2013 21st IEEE International Requirements Engineering Conference (RE)* (pp. 125-134). IEEE.
13. Guzman, E., Maalej, W. (2014). How do users like this feature? A fine-grained sentiment analysis of app reviews. In *2014 IEEE 22nd international requirements engineering conference (RE)* (pp. 153-162). IEEE.
14. Harman, M., Jia, Y., Zhang, Y. (2012). App store mining and analysis: MSR for app stores. In *2012 9th IEEE Working Conference on Mining Software Repositories (MSR)* (pp. 108-111). IEEE.
15. Haggag, O., Grundy, J., & Hoda, R. (2024). Towards enhancing mobile app reviews: A structured approach to user review entry, analysis and verification. In *Evaluation of Novel Approaches to Software Engineering (ENASE) Conference* (pp. 598-604).
16. Hoda, R. (2023). Technical briefing on socio-technical grounded theory for qualitative data analysis. In *2023 IEEE/ACM 45th International Conference on Software Engineering: Companion Proceedings (ICSE-Companion)* (pp. 344-345). IEEE.

17. Hoda, R. (2021). Socio-technical grounded theory for software engineering. *IEEE Transactions on Software Engineering*, 48(10), 3808-3832. IEEE.

18. Iacob, C.,  Harrison, R. (2013). Retrieving and analyzing mobile apps feature requests from online reviews. In *2013 10th working conference on mining software repositories (MSR)* (pp. 41-44). IEEE.

19. Mayzlin, D., Dover, Y.,  Chevalier, J. A. (2014). Promotional reviews: An empirical investigation of online review manipulation. *American Economic Review*, 104(8), 2421-2455.

20. Palomba, F., et al. (2017). User reviews matter! Tracking crowdsourced reviews to support evolution of successful apps. In *2017 IEEE/ACM 39th International Conference on Software Engineering (ICSE)* (pp. 3-13). IEEE.

21. Ott, M., et al. (2011). Finding deceptive opinion spam by any stretch of the imagination. In *Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies-Volume 1* (pp. 309-319).

22. Brown, T. B., et al. (2020). Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.

23. Liu, Y., et al. (2019). RoBERTa: A robustly optimized BERT pretraining approach. *arXiv preprint arXiv:1907.11692*.

24. Tian, Y., et al. (2019). Uber ATG: Data driven approaches to computing accuracy for urban computing applications. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 33, 1398-1405.

25. Medhat, W., Hassan, A.,  Korashy, H. (2014). Sentiment analysis algorithms and applications: A survey. *Ain Shams Engineering Journal*, 5(4), 1093-1113.

26. Charmaz, K. (2014).

27. Pagano, D.,  Maalej, W. (2013). User feedback in the app store: An empirical study. In *2013 21st IEEE International Requirements Engineering Conference (RE)* (pp. 125-134). IEEE.

28. Iacob, C.,  Harrison, R. (2013). Retrieving and analyzing mobile apps feature requests from online reviews. In *2013 10th working conference on mining software repositories (MSR)* (pp. 41-44). IEEE.

29. Guzman, E.,  Maalej, W. (2014). How do users like this feature? A fine-grained sentiment analysis of app reviews. In *2014 IEEE 22nd international requirements engineering conference (RE)* (pp. 153-162). IEEE.

30. Hu, N., Pavlou, P. A.,  Zhang, J. (2012). On the value of online product reviews: A theoretical model and empirical analysis. *MIS Quarterly*, 36(1), 291-314.

31. Ott, M., et al. (2011). Finding deceptive opinion spam by any stretch of the imagination. In *Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies-Volume 1* (pp. 309-319).

32. Mayzlin, D., Dover, Y.,  Chevalier, J. A. (2014). Promotional reviews: An empirical investigation of online review manipulation. *American Economic Review*, 104(8), 2421-2455.

33. Palomba, F., et al. (2017). User reviews matter! Tracking crowdsourced reviews to support evolution of successful apps. In *2017 IEEE/ACM 39th International Conference on Software Engineering (ICSE)* (pp. 3-13). IEEE.

34. Brown, T. B., et al. (2020). Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.

35. Liu, Y., et al. (2019). RoBERTa: A robustly optimized BERT pretraining approach. *arXiv preprint arXiv:1907.11692*.

36. Medhat, W., Hassan, A.,  Korashy, H. (2014). Sentiment analysis algorithms and applications: A survey. *Ain Shams Engineering Journal*, 5(4), 1093-1113.

37. Charmaz, K. (2014). *Constructing grounded theory*. Sage.
38. Corbin, J.,  Strauss, A. (2015). *Basics of qualitative research: Techniques and procedures for developing grounded theory*. Sage publications.