



Classification of quality characteristics in online user feedback using linguistic analysis, crowdsourcing and LLMs[☆]

Eduard C. Groen^{a,b,*,*}, Fabiano Dalpiaz^b, Martijn van Vliet^b, Boris Winter^b, Joerg Doerr^{a,c}, Sjaak Brinkkemper^b

^a Fraunhofer Institute for Experimental Software Engineering IESE, Fraunhofer-Platz 1, 67663 Kaiserslautern, Germany

^b Utrecht University, Department of Information and Computing Sciences, Princetonplein 5, 3584 CC Utrecht, The Netherlands

^c University of Kaiserslautern–Landau, Gottlieb-Daimler-Straße, 67663 Kaiserslautern, Germany

ARTICLE INFO

Dataset link: <https://doi.org/10.5281/zenodo.15604749>

Keywords:

Crowd-based RE
Crowdsourcing
Large language models
Online user reviews
Quality requirements
Requirements engineering
User feedback analysis

ABSTRACT

Software qualities such as *usability* or *reliability* are among the strongest determinants of mobile app user satisfaction and constitute a significant portion of online user feedback on software products, making it a valuable source of quality-related feedback to guide the development process. The abundance of online user feedback warrants the automated identification of quality characteristics, but the online user feedback's heterogeneity and the lack of appropriate training corpora limit the applicability of supervised machine learning. We therefore investigate the viability of three approaches that could be effective in *low-data* settings: language patterns (LPs) based on quality-related keywords, instructions for crowdsourced micro-tasks, and large language model (LLM) prompts. We determined the feasibility of each approach and then compared their accuracy. For the complex multiclass classification of quality characteristics, the LP-based approach achieved a varied precision (0.38–0.92) depending on the quality characteristic, and low recall; crowdsourcing achieved the best average accuracy in two consecutive phases (0.63, 0.72), which could be matched by the best-performing LLM condition (0.66) and a prediction based on the LLMs' majority vote (0.68). Our findings show that in this low-data setting, the two approaches that use crowdsourcing or LLMs instead of involving experts achieved accurate classifications, while the LP-based approach had only limited potential. The promise of crowdsourcing and LLMs in this context might even extend to building training corpora.

1. Introduction

Engaging end-users and capturing their requirements are crucial for a software system's success (Bano et al., 2017). A popular source of information is online user feedback (Astegher et al., 2023), which can be analyzed for user statements that express or pertain to requirements (Dąbrowski et al., 2022). Approaches to (semi)automatically analyzing online user feedback in requirements engineering (RE) are commonly referred to as CrowdRE (Groen et al., 2017b).

In addition to many CrowdRE works on functional requirements (Khan et al., 2019), some researchers have focused on identifying quality requirements (e.g., Groen et al., 2017a; Jha and Mahmoud, 2017; Lu and Liang, 2017). Online user feedback addressing software product qualities can inform about *how well* the software delivers its functions, i.e., a system's *qualities*. In particular, online user feedback

has been found to address quality characteristics far more than functional aspects (Groen et al., 2017a). Especially negative experiences and opinions—such as poor usability or instability—have an impact on user satisfaction (Ceaparu et al., 2004; Groen et al., 2017a; Hertzum and Hornbæk, 2023).

However, any approach aimed at classifying online user feedback—from manual to fully automated ones—faces difficulties due to informal language and finer nuances (e.g., Internet slang, expressions, sarcasm, emojis), author characteristics (e.g., laypeople, non-native English speakers and multilinguality, fakes), and general poor writing (e.g., sloppy use of punctuation, poor spelling, abundance of typos; Groen et al., 2018; Williams and Mahmoud, 2017). Moreover, users provide a personal account in prose, the relevant software quality that causes (dis)satisfaction is often implicit, and users are not always able to describe problems well.

[☆] Editor: Neil Ernst.

^{*} Corresponding author at: Fraunhofer Institute for Experimental Software Engineering IESE, Fraunhofer-Platz 1, 67663 Kaiserslautern, Germany.

E-mail addresses: eduard.groen@iese.fraunhofer.de (E.C. Groen), f.dalpiaz@uu.nl (F. Dalpiaz), m.vanvliet@uu.nl (M.v. Vliet), boris@prostrive.io (B. Winter), joerg.doerr@iese.fraunhofer.de (J. Doerr), s.brinkkemper@uu.nl (S. Brinkkemper).

Research problem. Classifying online user feedback along quality dimensions faces unique challenges. This is due not only to the ambiguity of the data to be classified but also to the way in which taxonomies model quality dimensions. In computer science, the ISO/IEC 25010 standard (ISO, 2011) is the de facto standard with respect to software product qualities. It provides a taxonomy of eight characteristics—*functional suitability*, *performance efficiency*, *compatibility*, *usability*, *reliability*, *security*, *maintainability*, and *portability*—each with two or more subcharacteristics.¹ Quality characteristics of software are typically specified in the form of *quality requirements*,² which are crucial for the success of software products. In RE, the ISO 25010 taxonomy is commonly used to organize quality requirements in a specification (Pohl and Rupp, 2015; Glinz et al., 2023), which is why we base our work on this particular taxonomy. However, because laypeople writing online user feedback often do not describe the circumstances in sufficient detail, it can be difficult to differentiate between these quality characteristics. For example, data getting corrupted might be primarily considered as *integrity* (a subcharacteristic of *security*), but if this happened during or because of an update, it could also be considered *replaceability* (a subcharacteristic of *portability*). The more ambiguous the feedback, the harder it is to correctly identify the affected quality. For instance, a statement such as “app doesn’t work” could imply a problem with *installability*, *reliability*, *compatibility*, *interoperability*, or *security*, or maybe even reflect the user failing to understand how the app is operated.

Although manual annotation by domain experts is the most accurate way of performing the difficult task of identifying software qualities in online user feedback, it is cognitively strenuous, time-consuming, and costly. Especially as the number of user reviews grows, automated analysis becomes indispensable (Groen et al., 2018). Finding a reliable alternative to expert-based annotations is an ongoing concern in CrowdRE. There are only a limited number of smaller annotated datasets (for reviews, see Dąbrowski et al., 2022; Reddy Mekala et al., 2021), and these were created based on separate tagging schemata. This poses limitations to the use of machine learning (ML) classifiers, whose performance degrades when applied to unseen data, due to the heterogeneity of RE data (Dell’Anna et al., 2023). Hence, we call this domain a *low-data setting*, referring to the absence of large and reliably labeled training corpora that can be used for training traditional ML and DL algorithms. In this setting, and also based on initial experimentation, we decided to explore alternatives that do not require training data.

Solution approach. To address the research challenge, we investigated the viability of approaches that do not require substantial training data. Specifically, we tested each approach on the task of classifying online user reviews from app stores. We selected three heterogeneous low-data approaches for this purpose:

- **Language patterns (LPs).** Inspired by related work on linguistic approaches (see Section 6.2), we adapted a method that was designed for eliciting non-functional requirements (NFRs), called the *NFR Method* (Dörr, 2011), and constructed an approach for classifying software qualities based on predefined (combinations of) keywords without machine learning. We elicited taxonomies of keywords & phrases through structured expert workshops, and then queried these over the online user feedback through LPs based on regular expressions. We considered this approach as a baseline against which to test the other approaches, which is important for drawing meaningful inferences.

¹ The experimental data on which this work reports was constructed when ISO/IEC 25010:2011 was in effect. By now, it has been superseded by the ISO/IEC 25010:2023 standard (ISO/IEC, 2023), which introduces *safety* as a ninth quality characteristic, renames *usability* to *interaction capability*, and *portability* to *flexibility*.

² Glinz (2007) has suggested *non-functional requirements* (NFRs) to be the overarching term for quality requirements and constraints. Our work does not address constraints, so we will speak of quality requirements where appropriate.

Table 1

Overview of the ISO 25010 quality characteristics and their mapping onto the classes used for answering RQ2 and RQ3.

ISO 25010 quality characteristic	Kyōryoku quality aspect	Short name
Compatibility & Portability	System support feedback	Compatibility
Usability	User-friendliness feedback	User-friendliness
Security	Security feedback	Security
Performance efficiency	Performance feedback	Performance
Reliability	Stability feedback	Stability
–	Feature request	Feature
–	None of the above/Other	None

- **Crowdsourcing.** We extended our previous research on the *Kyōryoku* method, presented in van Vliet et al. (2020), to assess how accurately a crowd of paid laypeople can classify online user feedback. The *Kyōryoku* method was specifically designed to cater to crowd workers without prior knowledge of RE, who are, among other things, not familiar with the distinction between features and the qualities of these features. The crowd workers are given a brief training and subsequently perform a micro-task. The instructions include examples, in line with the concept of few-shot learning in ML.
- **Large language models (LLMs).** We developed a pipeline to prompt an LLM with context information and instructions to perform a classification task. Using this approach, we investigated how accurately different conditions based on LLM, prompt type, and learning strategy can classify online user feedback.

Our categories of software product *qualities* were based on ISO 25010. In some cases, they were renamed in accordance with the taxonomy of Glinz (2007) to reduce task complexity and maximize comprehensibility for laypeople (cf. van Vliet et al., 2020). To help crowd workers understand better how *reliability* is distinct from *performance*, we chose *stability*, which reflects the most frequently addressed aspect of *reliability* in online user feedback (Groen et al., 2017a). Because sentences in online user feedback often describe whether or not the software supports a particular device or platform, it is difficult to determine whether this pertains to *compatibility* or *portability*, so we combined them into a single class. Two ISO 25010 characteristics are not explicitly included: *Functional suitability* is difficult to comprehend and is implicitly covered by *feature request*, and *maintainability* is not visible to the user, but indirectly affects other qualities (Groen et al., 2017a). The category *feature request* is included, so the two main requirements types—functional and non-functional—are considered in the classification (cf. Glinz, 2007).

Main research question. We investigated each of the three approaches in an in-lab evaluation study in which we sought to achieve an optimal configuration for each to answer this work’s main research question (MRQ):

MRQ. Which low-data approach is most accurate in identifying quality aspects in online user feedback?

Contributions. The goal of this paper is to empirically investigate the effectiveness of and trade-offs between three low-data classification approaches. In doing so, this paper makes the following contributions:

- We derive a keyword meta-model to help identify statements about quality aspects in online user feedback.
- We extend the crowdsourced user feedback classification method *Kyōryoku* (van Vliet et al., 2020) with classifications into further quality aspects.
- We provide a state-of-the-art LLM pipeline to classify online user feedback according to requirements relevance and into particular qualities.
- We compare the three approaches that were fine-tuned to identify quality-related information in online user feedback.

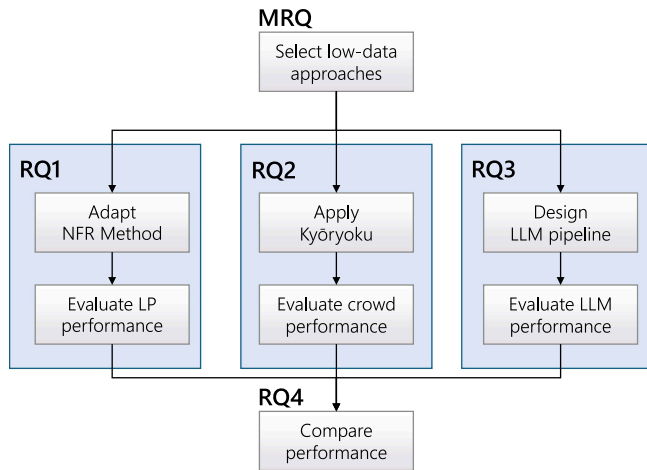


Fig. 1. Illustration of the relationship between the main research question (MRQ) and the research questions (RQs).

Paper outline. We begin by outlining our research methodology in Section 2. We then present the results for each research question in Section 3 and discuss our findings in Section 4. In Section 5, we present a discussion of the main threats to validity. Section 6 provides an overview of the relevant literature, and in Section 7, we conclude. An online appendix provides supplementary material including additional content, experimental artifacts, and spreadsheets with data and results (Groen et al., 2025).

2. Research methodology

In this section, we first detail the research questions (RQs) into which we decomposed our MRQ and the methods we used (Section 2.1); then we introduce the dataset and gold standards (Section 2.2) and describe the metrics and data analysis techniques we applied (Section 2.3). After that, we describe the experimental configuration for each of our four RQs in Sections 2.4–2.7.

2.1. Research questions & methods

In our investigation of which low-data approach is most accurate in identifying quality aspects in online user feedback, we selected and developed three candidate approaches based on the literature. To achieve this, we followed the design science approach for each approach, which Wieringa (2014) distinguishes into two activities: a *design activity* in which an artifact or other object is improved upon to address a practical problem, and an *investigative activity* to investigate and scientifically evaluate this artifact in context. Our comparison takes into account that each of the chosen approaches has particular benefits and opportunities for research or practice, but also comes with its unique challenges and drawbacks. Therefore, we decomposed our MRQ into four RQs, which address the design-investigate cycles for each approach individually (RQ1–RQ3) before comparing them (RQ4), as visualized in Fig. 1.

RQ1. How to adapt the NFR Method in order to derive LPs for the identification of quality characteristics in online user feedback?

We adapted the *NFR Method* to derive LPs for the identification of quality characteristics in online user feedback, which is an empirical approach that aims to ensure that all quality requirements are elicited and correctly specified (Dörr, 2011). This is a difficult challenge in industry (Chung and do Prado Leite, 2009), which it addresses by incorporating concepts of various NFR approaches, including *Goal-Oriented*

RE (van Lamsweerde, 2001) and the *NFR Framework* (Chung et al., 2012). The process adapts the backlog of (reusable) experience-based artifacts such as high-level quality attributes, custom quality models, templates, and checklists to a specific business and domain context. The process consists of two phases, each encompassing several activities. In Phase 1, elicitation and specification are prepared by prioritizing the quality characteristics, preparing models, identifying potential conflicts between qualities, and deriving checklists and templates. In Phase 2, the quality requirements are elicited and dependencies between them are identified. Our online appendix (Groen et al., 2025) presents a detailed description of the NFR Method and its process.

In our work, we operationalized the NFR Method’s strategy for the context of online user feedback about mobile apps. We first elicited relevant textual building blocks denoted as *keywords and phrases* (K&Ps) and then refined them into *language patterns* (LPs), which we used as queries to identify associated statements in online user feedback. In order to assess the merit of the designed approach based on LPs, we introduce and address research sub-question **RQ1a: How effective are language patterns based on quality-related keywords in identifying quality characteristics?**

RQ2. How to design micro-tasks so that lay workers recruited through crowdsourcing can achieve good results classifying quality aspects in online user feedback?

To investigate the potential of lay workers to classify quality aspects in online user feedback, we further researched the Kyōryoku³ method introduced in van Vliet et al. (2020), which describes a crowdsourced annotation process to identify requirements-related content in online user reviews about software apps. Following CrowdForge, a general-purpose framework on crowd work (Kittur et al., 2011), Kyōryoku involves a stepwise classification process through a series of micro-tasks. In crowdsourced work, *micro-tasks* are structured and simplified data extraction decision workflows that are performed by *crowd workers* in return for relatively small rewards (Retelny et al., 2014; Valentine et al., 2017). Because each micro-task is more granular than the preceding one, they become increasingly demanding and complex for crowd workers (cf. Schenk and Guittard, 2011; Gilardi et al., 2023), but this can be mitigated by making sure the micro-tasks’ instructions have been formulated to be simple enough for laypeople without relevant prior knowledge to understand them (Kittur et al., 2011). The method results in an annotated dataset that can, for example, be used as a gold standard to train ML classifiers. Each micro-task in Kyōryoku begins with a *job description*, followed by an *eligibility test* that needs to be passed in order to proceed to the actual *annotation task*.⁴ In order to investigate the performance of crowd workers in different configurations, we proposed two research sub-questions that allow us to evaluate two main design choices of our approach:

RQ2a. How does decomposing a classification task into smaller, consecutive classification tasks affect the performance of crowd workers?

This sub-question compares the performance of the crowd in a condition in which the crowd workers assigned all quality aspects to a condition consisting of two micro-tasks in sequence.

RQ2b. Can crowd workers correctly differentiate between all quality aspects?

This sub-question analyzes whether differences exist in the ability to recognize the various quality aspects.

³ Kyōryoku (協力) is a Japanese term for *collaboration*: Literally, it combines *strength* (力) with *cooperation* (協).

⁴ The job descriptions, test questions, and data sample are available in the online appendix of van Vliet et al. (2020) at doi:10.5281/zenodo.3754721.

RQ3. How to design an LLM pipeline so that it can achieve good results in classifying quality aspects in online user feedback with little training data?

To determine whether an LLM pipeline is capable of correctly classifying online user feedback into quality dimensions with little training data (RQ3), we performed an experiment with a series of classification tasks similar to RQ2, in which we varied the *large language model* (LLM), the prompt type, and the learning strategy. Our research design motivates two research sub-questions for RQ3:

RQ3a. How does increasing task complexity affect the performance of LLMs?

This sub-question compares the performance of LLMs in phases with increasing complexity: P1, P2, and P3'.

RQ3b. How does the use of engineered prompts and examples affect the performance of LLMs?

This sub-question assesses the impact of the experimental factors prompt type and learning strategy.

RQ4. Which low-data approach is most suited for classifying quality aspects in online user feedback, considering effectiveness and trade-offs?

We compared the three low-data classification approaches in terms of their trade-offs and effectiveness in classifying online user feedback into quality aspects to help us determine which approach is most accurate. We broke RQ4 down into two research sub-questions:

RQ4a. How suitable is each approach for determining requirements relevance?

This sub-question compares the performance of the crowd to that of LLMs in terms of filtering out online user feedback irrelevant to RE (P1 & P2).

RQ4b. How suitable is each approach for distinguishing between the quality aspects?

This sub-question compares the performance of all three approaches on the task of assigning the correct quality characteristics (P3 & P4).

2.2. Dataset & gold standards

In order to allow comparisons between the three approaches considered, we used a shared subset of the large dataset first introduced in Section II.A in Groen et al. (2017a). This dataset contains 132,194 online user reviews curated from six *app categories*—the five categories found by Pagano and Maalej (2013) to attract the most reviews, plus the emerging category “smart products”—with two *apps* per category—one paid, one free—and from three international English-language *app stores*—Apple App Store, Google Play, and Amazon.com—for a total of 36 sources. Our subset omits Amazon’s defunct app store for replicability reasons, and “smart products” apps because these were found to substantially differ from the other categories (Groen et al., 2017a), retaining 122,899 user reviews from 20 sources.

We originally developed our gold standards in van Vliet et al. (2020), for which we took a data sample of 1000 reviews that was systematically stratified across apps and app stores. This number was chosen to fit the job size limit of the account type used on the platform that we used in our study to answer RQ2. We created three gold standards for a multi-phase analysis. We took the dataset from van Vliet et al. (2020), but revised it to make sure that each classification was judged by three raters, for which we reconciled disagreements:

- The gold standard for Phase P1 distinguishes user reviews into helpful and useless ones from the perspective of the developer of an app.
- The gold standard for Phase P2 distinguishes helpful user reviews from P1 into relevant and useless sentences from the perspective of the developer of an app.

- The gold standard for Phases P3 & P4 distinguishes helpful sentences from P2 into the quality aspects listed in Table 1.

Table 2 shows how three online user reviews perpetuated through the gold standards. The 1000 reviews used for the gold standards and which served as our *test set* were omitted from the dataset to prevent overfitting. All artifacts related to the data and the gold standards are available from van Vliet et al. (2020) and in the online appendix to this paper (Groen et al., 2025).

After identifying an inconsistency in the data sample of van Vliet et al. (2020), in which text was erroneously inserted into some items, we performed data sanitation by manually removing reviews or sentences that could have been classified differently by the crowd. This resulted in the omission of 23 reviews in P1 (0.02%), 635 sentences in P2 (51.12%), 112 in P3 and P3' (16.40%), and 247 (19.03%) in P4. This appears to have only affected our ability to draw conclusions about the *useless* classifications in P2; see our online appendix (Groen et al., 2025) for details on the data sanitation process and tables comparing the data.

2.3. Metrics & data analysis

For each approach, we compared the classifications against the gold standards to determine the true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN). For each phase in RQ2, we created a combined prediction of the crowd workers’ judgments by majority vote, with a tie constituting a multi-label decision. Inspired by the recommendation of Mizrahi et al. (2024) to evaluate LLMs by averaging them across multiple prompts, we similarly calculated a prediction by majority vote for RQ3. Based on these decision scores, we gauged the crowd’s performance using the three primary performance metrics for evaluating classifier performance on a classification problem (cf. Dell’Anna et al., 2023):

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

Of these, *accuracy* describes the fraction of the entire space that is classified correctly by a classifier, whether as a true positive or as a true negative; *precision* measures the degree to which only the relevant answers are found; and *recall* measures the degree to which all the relevant answers are found (Berry, 2021). Because our work explores the effects of both precision and recall, we report these scores instead of F_1/F_β harmonized means of these two metrics. To compare the accuracy of the classification against the gold standards, we calculated correctness scores, constructed confusion matrices to identify where the most errors in a particular condition were made, and which two classes were confused most often. For cases in which the gold standard of P3 & P4 allowed for multiple correct responses, the correct predictions were counted. If the majority vote was tied, resulting in a multi-label output, this response was considered correct if at least one prediction matched the gold standard; if all predictions were incorrect, these were counted as separate penalties for each incorrect class. We provide the macro average for precision and recall, and the mean score representing the accuracy value. To visualize differences in classifier performance, we created receiver operating characteristic plots (ROC plots; Fawcett, 2006), which contrast *sensitivity*—which is synonymous with recall—with *specificity*:

$$Specificity = \frac{TN}{TN + FP}$$

Although for statistically analyzing classifiers—usually ML algorithms—Dell’Anna et al. (2023) suggest repeated-measures designs, we consider the configuration of the LLM conditions for RQ3 to be a between-subjects design because we did not control for non-determinism by

Table 2

Examples of online user reviews in the gold standards, with classifications in square brackets. Items classified as useless are not perpetuated to the next phase. The sentence splitter for Phase P2 splits sentences based on punctuation. In Phase P2, sentences that rely on the context to be understood are judged useless. In Phase P3, multiple classifications were only assigned if a sentence was interpreted to ascribe about equal value to two different aspects.

Gold standard P1	Gold standard P2	Gold standard P3 & P4
Can't access settings or feedback What good is an application if you can't change the settings or send feedback through it? I have the same problem with office lens. Typical arrogant Microsoft. [Helpful]	Can't access settings or feedback What good is an application if you can't change the settings or send feedback through it? [Helpful] I have the same problem with office lens. [Useless] Typical arrogant Microsoft. [Useless]	Can't access settings or feedback What good is an application if you can't change the settings or send feedback through it? [Feature request, User-friendliness] – –
App addict. This is the one and only and best app ever! Can't live without it! [Useless]	–	–
Crashes when I open it & terrible lag I've used this app for over 2 yrs & have loved it until now. Every time I try to reply to someone the app closes. Also, the lag is terrible. This has been the best app until lately. [Helpful]	Crashes when I open it & terrible lag I've used this app for over 2 yrs & have loved it until now. [Helpful] Every time I try to reply to someone the app closes. [Helpful] Also, the lag is terrible. [Helpful] This has been the best app until lately. [Useless]	Crashes when I open it & terrible lag I've used this app for over 2 yrs & have loved it until now. [Stability] Every time I try to reply to someone the app closes. [Stability] Also, the lag is terrible. [Performance] –

altering the configuration of hyper-parameters—which is not possible for ChatGPT—or by sampling multiple outputs, and the various experimental factors inherently intended the LLMs to act differently across conditions.

Our data was not normally distributed; Mardia's test for multivariate normality (Mardia, 1970) revealed significant skewness from symmetry, $b_{1,p} = 1290.01$ (P1), 1534.57 (P2), 180.29 (P3'), $p < .001$, and significant kurtosis implying shorter tails than expected, $b_{2,p} = -44.81$ (P1), -53.39 (P2), -39.64 (P3'), $p < .001$. In order to measure the main effects and the interaction effects of the three factors (prompt type, learning strategy, and LLM) on classification accuracy for Phases P1, P2, and P3' and to compare LLM accuracy by phase, we used *binomial Generalized Linear Models* (GLM; Hosmer et al., 2013), which extend traditional linear regression to binary data for which normality cannot be assumed and allow the modeling of non-linear relationships to estimate the effect of our experimental factors. Because we are also interested in how likely it is that a change in a condition will alter the outcome, we quantified the strength of the effects using the *Odds Ratio* (OR), which exponentially transforms the GLM's regression coefficient β through the logistic function $OR = e^{\beta}$, and which serves as an indicator of the multiplication factor by which a change in condition is expected to increase or decrease (Kutner et al., 2005). Because we found some results that would have been significant had the data been normalized, we also report results that show a trend toward significance. Because almost from the beginning, we had indications that the LLM used had the greatest impact on the results, we often contrasted the performance of ChatGPT 4 Legacy and ChatGPT 4o.

To answer RQ4, we compared the performance of the three approaches in several ways. To compare classifier performance, we contrasted the means with Tukey's Honestly Significant Difference (HSD) test for pairwise comparisons (Tukey, 1949). To compare the performance of the classifiers by class, we analyzed them individually as binary predictions. This did impact the accuracy of the classifiers based on crowdsourcing and ChatGPT because it entailed assigning penalties for each missed multi-label option in the gold standard and each incorrectly provided multi-label prediction. Because a comparison of performance on individual classes constitutes a repeated-measures design, we compared the performance of the classifiers on individual classes using the Friedman Test (Dell'Anna et al., 2023). For significant χ^2 scores, we further explored this using Wilcoxon signed-rank tests for pairwise comparisons with Bonferroni-adjusted p -values (cf. Demšar, 2006). We contrasted the classifier with the highest mean per class against the other scores using binomial GLMs.

2.4. Experimental configuration for RQ1 — Language patterns

To answer RQ1, we elicited keywords and phrases (K&Ps) through a series of workshops adapted from the NFR Method (Dörri, 2011), based

on which we created language patterns (LPs). To systematically elicit K&Ps for the quality characteristics of ISO 25010, we conducted six elicitation workshops with a maximum duration of 1.5 h with a total of 24 participants (11 female, 13 male) solicited from among Computer Science students at two universities, doctoral candidates in Computer Science, and scientific staff at Fraunhofer IESE, a research institute. The workshop addressing the characteristics of *usability*, *security*, and *maintainability* was attended by the same experts. Most of the participants were German (11) or Dutch (7) and had good self-reported English proficiency. Despite the fact that all of the participants had been using a smartphone for many years, only five had ever written a user review. Due to the COVID-19 pandemic, the final two workshops were held remotely. Each workshop consisted of five steps: (1) briefing and formalities, (2) aligning on the quality characteristic, (3) establishing the quality model, (4) defining K&Ps, and (5) debriefing. To this end, the participants received the ISO 25010 definitions of their workshop's quality characteristic and its subcharacteristics, written guidelines designed to ensure lively interaction during the workshop, and templates to develop participants' understanding of the quality characteristic and to collect the K&Ps.

After the workshops, we documented the outcomes and calculated descriptive statistics. We recoded the elicited K&Ps into LPs using regular expression notation, so that the K&Ps can be searched automatically through queries. Because the K&P were obtained during brainstorming sessions in which the participants had little time to reflect on the results, and because the participants looked at only one characteristic, two authors served as judges to realign the classification, which included the inclusion of synonyms, variations, negated antonyms, and forbidden words to prevent false positives (FPs). This resulted in 248 LPs, which we queried over our analysis set (Round 1). We analyzed all results for the LPs with up to 100 matched statements, or a stratified random sample of 100 matched statements for LPs with more matches, so that in total, we reviewed 7467 statements. We changed the regular expression syntax to exclude FPs, where possible, to increase precision. An example of an LP for the *portability* sub-characteristic *replaceability* is: `(?i)(?!should|could|would|th)(is|are|has|have)(a|an|l)(far|much|)(more|)(improv|upgrade|faster|quicker)`. We repeated the query and the analysis over a sample of 6149 statements (Round 2). LPs with a precision below 0.50 were discarded. The final set of 242 LPs was used to query our test set to measure precision and recall against the gold standard for P3 & P4. The online appendix (Groen et al., 2025) provides a detailed description of the relatively complex LP encoding process and a primer on regular expression notation.

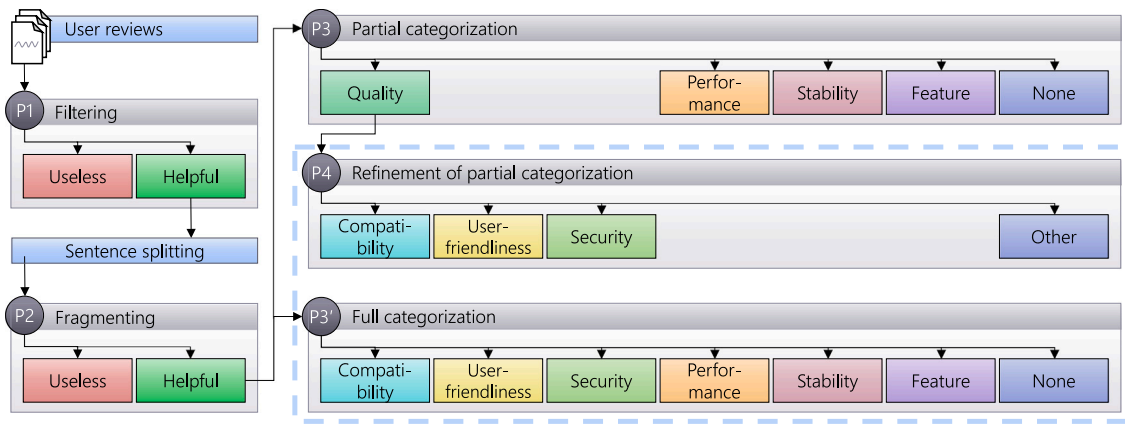


Fig. 2. Overview of the classifications in the various phases performed using the Kyōryoku method. Our work extends (van Vliet et al., 2020) with two phases, P4 and P3', shown with a dashed line.

2.5. Experimental configuration for RQ2 — Crowdsourcing

To answer RQ2, we conducted a single-group experiment for which we recruited crowd workers through the online crowdsourcing marketplace *Figure Eight*,⁵ which allows crowd workers to perform jobs by assigning micro-tasks in exchange for fixed-price monetary rewards. This section presents an abridged overview of the method presented in van Vliet et al. (2020), extended by Phases P4 and P3'.

We selected *Figure Eight* because of its support for data categorization tasks and its many embedded quality control mechanisms. Following the Kyōryoku method, a set of 1000 user reviews was classified in several consecutive micro-tasks, initially distinguishing between requirements-relevant and requirements-irrelevant (cf. Sihag et al., 2023), and later between various quality characteristics. Fig. 2 visualizes the sequence of the phases, which are as follows:

- P1** The annotation task of P1 served to determine whether the body text of unprocessed user reviews could be *helpful* to software developers or should be considered *useless*.
- P2** The *helpful* reviews from P1 were split into individual sentences that the crowd workers in P2 distinguished into *helpful* and *useless*.
- P3** The *helpful* sentences obtained from P2 underwent a more fine-grained classification (see Table 1 in Section 1). In van Vliet et al. (2020), we reported on Phase 3 (P3), in which the crowd workers performed a partial classification in which the classes *compatibility*, *user-friendliness*, and *security* were combined into the class *quality feedback*.
- P4** In this work, the *quality feedback* class served as input for P4. The class *none of the above* was used for aspects such as general criticism and praise; in P4, this category was called *other* because based on the classification of the preceding P3, it already described a quality of some sort.
- P3'** We furthermore investigated the crowd workers' performance when classifying all classes at once in Phase 3 Primed (P3'), where the three quality aspects combined in P3 as *quality feedback* were presented together with the other quality aspects.

We refer to the first two phases, which focus on requirements relevance, as **P1 & P2**. Because P3 and P4 together constitute a classification of all quality aspects in two consecutive phases, we jointly refer to those as **P3→P4**. We collectively refer to the combination of Phases P3→P4 and P3' as **P3 & P4**. Table 3 shows the configuration for these five phases, with P3 & P4 enabling us to compare two variants: a single session with all classes (P3') and two consecutive sessions

(P3→P4). In *Figure Eight*, we opted for an open crowd selection policy; it is realistic to expect non-native English-speaking crowd workers to be capable of performing such a task, and their language proficiency could be ensured through the eligibility test. Three crowd worker judgments were sufficient to determine the majority vote for the binary classifications in P1 and P2. For the multiclass classifications in P3 & P4, we based the majority vote on six crowd worker judgments. The total cost of our experiment was \$619.32, which includes *Figure Eight*'s 20% usage fee. The sessions were active for 10:44:47 h in total before reaching completion.

2.6. Experimental configuration for RQ3 — LLMs

To answer RQ3, we performed a between-subjects experiment with a 2 (prompt type) × 2 (learning strategy) × 2 (LLM) factorial design, for a total of eight conditions, comparing the following experimental factors:

1. Two *prompt types*: engineered (*Eng*) vs. Kyōryoku (*Kyō*);
2. Two *learning strategies*: few-shot (*Few*) vs. zero-shot (*Zer*);
3. Two *LLMs*: ChatGPT 4 Legacy⁶ (4) vs. ChatGPT 4o (4o).

Each condition was applied to the classification problems of Phases P1 (binary), P2 (binary), and P3' (multiclass, multi-label) described in Section 2.5. Due to the novelty of the domain, we compensated for the lack of established approaches with an elaborate pretest and by making improvements during the experiment. In the following, we will briefly discuss the most important decisions, which are further elaborated in our online appendix (Groen et al., 2025).

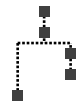
In a pretest, we compared Mixtral, Nous Hermes, ChatGPT 3.5, and the then-latest ChatGPT 4. We found that ChatGPT was the most capable of handling our prompts by their number of tokens and producing responses that were meaningful and correctly structured. Because this study serves to provide an estimation of the capability of LLMs rather than a comprehensive comparison of LLMs, we chose to perform this study using the state of the art in LLMs based on the literature available at the time (e.g., Gilardi et al., 2023). Our experiment was set up and trialed using ChatGPT 4 Legacy, which had been superseded by ChatGPT 4o by the time we conducted it. Because we had not studied the characteristics of ChatGPT 4o, which had just been released at the time, we decided to compare the performance of both LLMs. A downside is that the GPT models hosted by OpenAI are black boxes whose evolution cannot be controlled, which impairs their replication

⁶ Previously known as ChatGPT 4 Turbo, OpenAI renamed it ChatGPT 4 Legacy upon the release of its successor model, ChatGPT 4o.

⁵ Later acquired by Appen, <https://www.appen.com/>.

Table 3

Summary of the configuration for each phase offered on *Figure Eight*, along with runtime statistics per session. P1 and P2 consisted of two micro-tasks. The sequence representation illustrates how Phases P3→P4 are performed concurrently to Phase P3'.

Sequence	Phase	Launch (2019)	Judgments	Duration	Judgments per item	Contributors	Price per judgment (\$)	Total cost (\$)
	P1	May 7 & 15	50	1:48:40	3	123	0.03–0.04	130.80
	P2	May 23 & 29	50	1:10:12	3	235	0.02	106.32
	P3	June 13	50	2:19:21	6	145	0.02	117.60
	P4	Dec. 27	50	1:58:36	6	62	0.03	77.04
	P3'	Dec. 27	50	3:27:58	6	146	0.03	187.56
Total				10:44:47	711		619.32	

potential compared to other LLMs. We did not fine-tune or retrain models due to the unavailability of training data in our low-data setting. ChatGPT, on the other hand, usually hallucinates less than other LLMs, even though its hyper-parameter settings cannot be altered from the web interface. ChatGPT was configured without custom instructions or memory function. A detailed discussion of what implications the choice of ChatGPT had on the validity of our work is provided in Section 5.

Like Gilardi et al. (2023), who provided the exact same instructions to ChatGPT that they had given their crowd workers, we used the instructions from the Kyōryoku method that we used for our micro-tasks and our gold standards. However, because there was also some evidence that a prompt specifically engineered for ChatGPT might lead to better results (Brand, 2023; El-Hajjaji et al., 2024), we also engineered a prompt that included the examples that we also provided to the crowd workers, making this a few-shot prompting strategy. As an example, Fig. 3 shows the few-shot engineered prompt for P1. The literature suggests that the behavior of ChatGPT can be influenced both positively and negatively by the use of examples, so we also included a learning strategy factor in which we compared the effect of the prompt including examples (i.e., a few-shot prompting strategy) with one in which these were omitted (i.e., a zero-shot prompting strategy). We presented the input data using text-based prompts in batches of 100 items or however many remained: i.e., 10 batches for P1 (1000 reviews), 14 batches for P2 (1347 sentences), and 7 batches for P3' (625 sentences). The 31 batches across eight conditions amounted to a total of 248 treatments. To prevent contamination through learning effects, each treatment was performed in a new prompt session after restarting the Web browser. Because two responses per condition were found to often be sufficiently similar Gilardi et al. (e.g., 2023), we chose to collect one response per condition and instead focus on comparing other experimental factors that, to our knowledge, had not been empirically tested yet.

2.7. Experimental configuration for RQ4 — Comparison

To compare the performance of the three approaches (RQ4), we conducted a posttest-only nonequivalent groups quasi-experiment (Jhangiani et al., 2019) to obtain a performance benchmark between nonequivalent groups, with the classification problem as the independent variable. To predict requirements relevance in Phases P1 and P2, we compared the crowd's majority vote prediction (RQ2) against the best-performing ChatGPT condition and the majority vote prediction of the eight ChatGPT conditions (RQ3). Our LP approach (RQ3) did not predict requirements relevance. Because ChatGPT had an even number of conditions, 104 predictions in P1 and 99 predictions in P2 of these binary classification tasks resulted in a 4:4 tie, which we omitted from the analysis to prevent undue penalties. We compared the performance of all three approaches over the gold standard for P3 & P4 for five conditions:

- *Language Patterns*. The prediction through the LP-based analysis through 242 LPs obtained using the NFR Method.
- *Crowd P3→P4*. The majority vote prediction of six crowd worker judgments in Phase P3→P4 of the crowdsourced approach using Kyōryoku.

Table 4

Number of K&Ps by ISO 25010 quality characteristic, and number of LPs (N) and precision (P) achieved per round. Red and green indicate the highest and lowest precision value by class.

	Elicited K&Ps	LPs round 1 (N, P)	LPs round 2 (N, P)
Compatibility	42	41 0.75	41 0.91
Portability	23	52 0.59	27 0.88
Usability	40	36 0.57	24 0.78
Security	51	31 0.54	20 0.92
Perf. efficiency	54	57 0.74	37 0.88
Reliability	24	31 0.69	21 0.92
Total	234	248 0.65	242 0.87

- *Crowd P3'*. The majority vote prediction of six crowd worker judgments in Phase P3' of the crowdsourced approach using Kyōryoku.
- *ChatGPT Best*. The best-performing ChatGPT condition <Kyō, Zer, 4o> in P3' by absolute positives in the LLM-based approach.
- *ChatGPT Majority Vote*. The majority vote prediction of the eight ChatGPT conditions in P3' in the LLM-based approach.

3. Results

In this section, we present the results to answer each of our four research questions.

3.1. Results for RQ1 — Language patterns

Table 4 summarizes the quantitative results of the workshop outcomes and the LPs encoded from the K&Ps, which are further detailed in the online appendix (Groen et al., 2025). The workshops produced a good number of 234 K&Ps, ranging from 23 for *portability* to 54 for *performance efficiency*. After encoding the 248 LPs, they matched 143,044 statements in Round 1 and achieved an overall weighted micro-average precision of 65%, meaning that two-thirds of the statements matched by all LPs were true positives (TPs). After making adjustments, the 242 LPs in Round 2 matched 77,935 statements (-45.5%). Precision improved considerably for all characteristics, reaching a micro-average of 0.87 overall. In most cases, a minor change could increase a particular LP's precision in Round 2, and like Cleland-Huang et al. (2007b), we found that certain keywords were weak indicators by themselves (e.g., “available” for *availability*). Almost all LPs improved; 115 of them even reached perfect precision. The five LPs from Round 2 that still had a precision below 0.5 were removed, leading to a final set of 237 LPs.

Table 5 shows the confusion matrix for the performance of the LPs compared to the gold standard for P3 & P4. Compared to Round 2, the LPs had a lower but still fair precision in the 0.40–0.97 range, but achieved low recall in the 0.06–0.50 range. Because the pattern-matching only assigned positives (i.e., statements it positively matched), we considered all missed classifications as *none*, resulting in limited precision (0.31) but high recall (0.95) for this class because it missed many statements but also matched some statements labeled in the gold standard as *none*. *Reliability* achieved good results, with a precision

```

# Binary Classification Task for ChatGPT

## Task Description
Classify user reviews from mobile app stores (Google Play Store, Apple App Store) as either "helpful" or "useless". The objective is to filter out spam and irrelevant reviews, aiding developers in focusing on constructive feedback.

## Context
Developers depend on user reviews for actionable feedback. The task aims to pre-process these reviews, filtering out useless reviews that are irrelevant, whilst retaining helpful reviews.

## Input Data Source
The data to classify is provided at the end of this prompt under 'Reviews to Classify'.

## Class Definitions
- **Useless**: Reviews that are not helpful for developers. They may contain spam, be off-topic, express feelings without elaboration, lack genuineness, or contain jokes.
- **Helpful**: Reviews that offer specific feedback, report bugs, suggest improvements, comment on updates, or provide constructive criticism.

## Examples for Each Class

**Class 1: Useless**
- Example: "I Really Like This!"
- Example: "My kids love it. Thanks"

**Class 2: Helpful**
- Example: "Newest version crashes when opening"
- Example: "Buggy and unreliable. Does not work often. Signs me out regularly."

Note: The examples provided are real user reviews. They may not always follow conventional punctuation or stylistic norms, reflecting the authentic and varied nature of user-generated content.

## Instructions for Classification
1. Assess each review in the input data.
2. Classify each review as either **Helpful** or **Useless** based on the definitions and examples provided.

## Output Format
Present results in a numbered list providing three values on a single line:
- **Review**: The first three words from the input data, for reference.
- **Judgment**: Your assigned classification as Useless or Helpful.
- **Confidence**: A qualitative assessment of your confidence in the classification as High, Medium, or Low.
Provide only the values specified. Do not provide a justification for your judgment.

## Additional Considerations
- Begin the classification process by directly addressing all reviews below. Avoid trial runs or partial classifications; instead, apply the classification criteria to the entire dataset in a single, comprehensive review process.
- Strictly adhere to the provided class definitions when classifying reviews. Ensure that the classification as either 'useless' or 'helpful' is firmly based on the criteria specified, without leniency. This is crucial for maintaining the integrity and accuracy of the classification process.
- Consider each review in isolation. When classifying, if a review's content is unclear, ambiguous, or incomprehensible, use your best judgment. Classify a review as 'Useless' if its content is so unclear, ambiguous, or incomprehensible that it would not reasonably make sense to a human reader.
- Approach the classification with the understanding that examples serve as guidelines but may not cover all scenarios. Use judgment to classify reviews that may not fit neatly into one category.
- In cases where a review contains both helpful and useless elements, classify it as 'Helpful'. This approach ensures that potentially valuable feedback is not overlooked.

## Reviews to Classify
...

```

Fig. 3. Abridged engineered prompt for Phase P1 with examples, i.e., <Eng, Few>. See Figure 2 in van Vliet et al. (2020) for a comparison to the Kyōryoku prompt, i.e., <Kyō, Few>. We omitted some examples for the sake of brevity. See the online appendix (Groen et al., 2025) for the full set of prompts.

Table 5

Confusion matrix comparing the LP results to the gold standard for Phases P3 & P4. All aspects not matched by an LP were labeled none/missed.

LPs	Gold standard									
	Compatibility	User-Friendl.	Security	Performance	Stability	None/Missed	Mult.(corr.)	Mult.(false)	Precision	Recall
Interop./Port.	4	2	0	0	1	0	0	3	0.40	0.06
Usability	1	11	2	0	0	0	1	0	0.81	0.10
Security	2	4	4	0	1	0	1	0	0.42	0.24
Performance	2	1	0	4	0	5	1	0	0.47	0.13
Reliability	0	1	0	0	50	1	5	0	0.97	0.50
None/Missed	62	112	11	42	50	140	0	28	0.31	0.95
Mult.(corr.)	0	1	0	2	1	0	3			
Mult.(false)	0	0	0	0	0	1		0		
Macro average									0.56	0.33
Accuracy										0.41

of 0.97 and a recall of 0.50. The combination of *interoperability* and *portability* achieved the poorest results, with a precision of 0.40 and a recall of 0.06.

3.2. Results for RQ2 — Crowdsourcing

We will first describe the crowd that we assembled through *Figure Eight* and the job they performed (Section 3.2.1), and then report on

the accuracy of the crowdsourced work in terms of precision and recall by classification output (Section 3.2.2) and agreement (Section 3.2.3).

3.2.1. Crowd demographics & job statistics

We gathered a large worldwide crowd through multiple crowd work channels associated with *Figure Eight*. A total of 711 unique crowd workers commenced participation in one of the seven micro-tasks, 555 (78.06%) of whom contributed judgments beyond the eligibility test

Table 6

Number of pages completed by each contributor per session. That two participants were able to complete six pages in Phase P3 is presumably due to an anomaly in Figure Eight.

No. of pages	Participants (N, %)		P1	P2	P3	P4	P3'
<1	21	2.95	2	11	1	2	5
1	150	21.10	16	89	24	9	12
2	71	9.99	14	28	5	10	14
3	45	6.33	11	17	4	6	7
4	62	8.72	27	15	8	1	11
5	360	50.63	53	75	101	34	97
6	2	0.28			2		
Total	711	100.00	123	235	145	62	146

and passed the quality checks. These 555 workers can be considered contributors. Most of the crowd workers were from countries where parts of the population live in poverty, such as Venezuela (Posch et al., 2022), which alone accounted for 33.47% of all contributors.

We received a total of 16,428 annotations (24,307 including trial data), with an average of 35.72 classified items per contributor. This data includes trial data and actual judgments. A contributor who completed all 50 judgments generated 14 trial entries; ten from the first page served as an aptitude test, while on the next four pages, one in ten items was a test question per page. Table 6 shows that on the whole, half of the participants who started a session completed all five pages.

Our experimental setting resulted in the two scenarios illustrated earlier in Table 3. For the sequence P1→P2→P3', we attracted 504 crowd workers over a time span of 6:26:50 h, with a total cost of \$424.68. For the sequence P1→P2→P3→P4, we attracted 565 crowd workers (+61) over a time span of 7:16:49 h (+49:59 min), with a total cost of \$431.76 (+\$7.08). Note that these numbers are based on the same values for P1 and P2. Table 6 shows that in P1 and especially P2, more participants abandoned the session early on, while in P3 & P4, most of the contributors completed all pages of the session. This is also seen in the higher share of incorrectly answered trial questions in P1 & P2, and in Table 6 by the higher share of judgments achieved in P3 & P4.

3.2.2. Crowdsourced classifications

We found that the crowd's performance was significantly more accurate in the condition with two consecutive micro-tasks P3→P4 than in the singular classification task P3' ($\beta = 0.4235, SE = 0.12, z = 3.49, p < .001, OR = 1.53$), with the effect size suggesting that splitting P3 & P4 into the two consecutive subtasks P3→P4 is 1.5× more likely to achieve better results. In both P3 and P3', the crowd classified far fewer sentences as *none* than in the gold standard, with the combined score for P3→P4 being slightly above the gold standard. Although detrimental to precision, having fewer classifications of *none* is likely to support better recall because fewer quality-related sentences are discarded. The crowd workers in P3 were more likely to classify sentences as *feature* and *stability*, while crowd workers in P3' chose *performance* nearly twice as often as in P3. A notable similarity is that in P3, the participants classified 32.05% of the sentences as *quality*, which hardly differs from the share of 29.54% in P3' for the equivalent combination of *compatibility*, *user-friendliness*, and *security*.

The confusion matrix for P3 in Table 7 shows that the crowd achieved the best precision and recall on *feature*, and high recall on *stability*, with fairly good precision. The majority of the mistakes involved *none* getting classified as *quality* (17 instances) and vice versa (46 instances). Thus, it is likely that potentially irrelevant sentences belonging to the *none* class are included in the *quality* category that was used as input for P4, which might explain why the crowd workers in P4 classified a considerable share of items as *other*. The crowd workers in P3 also frequently classified *stability* and *feature* (and to a lesser degree *performance*) as *quality*, which can be explained by the latter being a more generic class that was more likely to be selected in case of doubt.

For *none*, precision was fairly good, but recall was worst for all tags, while precision was lowest for *performance*.

Table 8 shows that in P4, the crowd workers performed best on *compatibility*. Performance was mainly affected by the crowd workers confusing *user-friendliness* with *none* (22 instances) and vice versa (20 instances), in part due to misattributing sentences about connection quality. Sentences that the crowd workers in P3 failed to classify as *none* continued to persist in P4 and were then often classified as *user-friendliness*. Recall on *security* was drastically reduced because two out of five classifications were misattributed as *user-friendliness*. The poor performance on *security* does not paint a reliable picture because of the low number of associated sentences in the gold standard for P3 & P4.

The confusion matrix for P3' in Table 9 shows that in this phase, the crowd performed best on *stability* and *feature*. The most obvious deviation from the gold standard is the large number of sentences misattributed to *performance*, most of which should have been classified as *user-friendliness* (48 instances) or *none* (41 instances) according to the gold standard, while only 33 sentences were correctly classified as *performance*. Fig. 4 visualizes the distribution of the judgments assigned by the crowd workers in P3 & P4, which highlights that *performance* was overclassified. Several sentences that should have been classified as *none* were classified as *user-friendliness* (20 instances). The crowd tended to assign *security* less often compared to the gold standard, possibly because its low salience in the gold standard caused its prominence to be low among crowd workers.

3.2.3. Crowd worker analysis

Table 10 shows how often the crowd workers provided a correct response in P1 & P2 according to the gold standards, with the largest category in each phase being where all three crowd workers gave the correct response. Although we found that P2 was slightly more difficult than P1, the agreement accuracy does not appear to have suffered.

Table 11 shows the accuracy scores achieved in P3 & P4 based on the number of crowd workers agreeing on the same class. When all six contributors agreed, the classification matched very well with the gold standard, with the lowest accuracy achieved in P3 (84.72%). Typically, accuracy quickly decreased as agreement between crowd members decreased. The exception to this is that the accuracy for tags that two out of six crowd workers agreed on was higher than when three out of six agreed. This can be explained by a higher likelihood that two or three classes were chosen, of which at least one was the correct answer. The only occurrence of complete disagreement was observed in P3', where the crowd workers assigned every possible tag except for *security*. The sentence was "If not being able to search within your notes (including text in images like you can on the pc version) is disappointing, than [sic] the complete lack of any search feature is a disaster". The situation in which three out of six crowd workers agreed on one class occurred most often, which sometimes involved a tie between two classes. These were especially observed for the classification of *none* and *user-friendliness*, as well as *none* and *compatibility*.

Table 11 also shows that the overall accuracy of the crowd workers in P3 (63%) outweighed that of P3' (55%), with the best results in the smaller sequential task P4 (72%), suggesting that a multi-stage classification will lead to more accurate results. With higher agreement strongly increasing precision, the quality of the results is likely to scale with the number of classifications per sentence. Thus, a possible solution to achieving better results is not to limit the number of crowd workers per job, but to instead keep the judgment phase active until a certain number of contributors agree.

3.3. Results for RQ3 — LLMs

In this section, we first provide general statistics of the experiment with LLMs (Section 3.3.1), before reporting on the classification performance on P1 and P2 (Section 3.3.2), and on P3' (Section 3.3.3).

Table 7

Confusion matrix comparing the results of the crowdsourced task of Phase P3 to the gold standard.

Crowd	Gold standard							Precision	Recall
	Quality	Performance	Stability	Feature	None	Multiple (correct)	Multiple (false)		
Quality	103	16	5	0	46	8	0	0.63	0.58
Performance	11	16	1	3	5	1	1	0.46	0.56
Stability	24	2	65	2	12	6	2	0.63	0.88
Feature	21	0	1	49	5	6	1	0.67	0.89
None	17	1	1	1	38	0	0	0.66	0.42
Multiple(correct)	21	8	8	5	21	3			
Multiple(false)	18	1	2	1	13		0		
Macro average Accuracy								0.61	0.67
								0.63	

Table 8

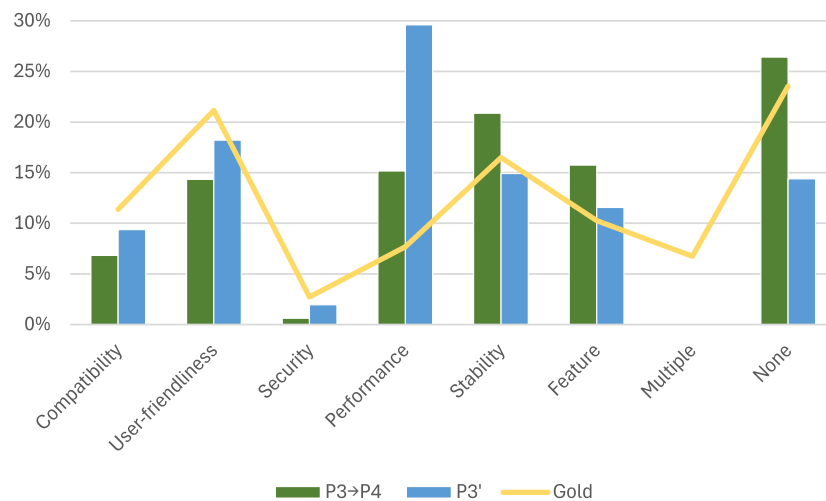
Confusion matrix comparing the results of the crowdsourced task of Phase P4 to the gold standard.

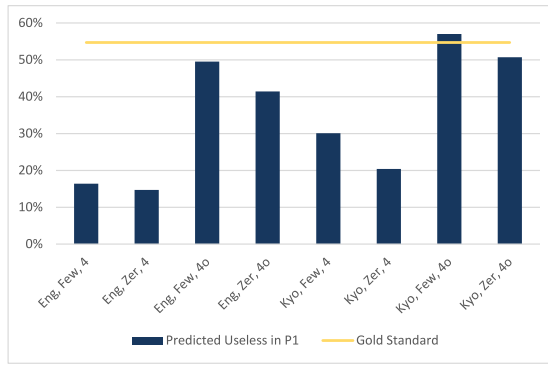
Crowd	Gold standard						Precision	Recall
	Compatibility	User-friendliness	Security	None	Multiple (correct)	Multiple (false)		
Compatibility	22	2	0	2	4	0	0.87	0.77
User-friendliness	5	40	2	22	3	0	0.60	0.70
Security	0	0	2	0	0	0	1.00	0.40
None	2	20	0	45	3	1	0.68	0.70
Multiple(correct)	1	11	0	11	1			
Multiple(false)	0	0	1	0		0		
Macro average Accuracy							0.79	0.64
							0.72	

Table 9

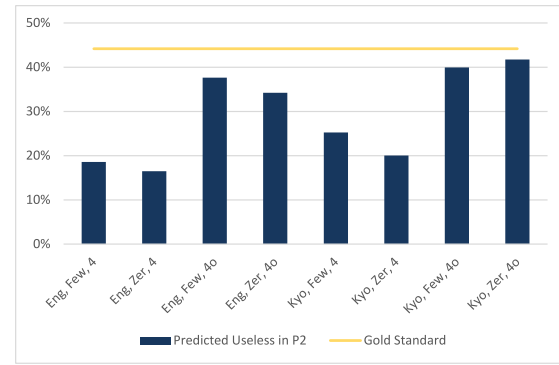
Confusion matrix comparing the results of the crowdsourced task of Phase P3' to the gold standard.

Crowd	Gold standard									Precision	Recall
	Compatibility	User-friendl.	Security	Performance	Stability	Feature	None	Mult.(corr.)	Mult.(false)		
Compatibility	19	0	2	2	2	1	7	2	0	0.58	0.44
User-friendl.	7	34	2	1	1	3	20	1	1	0.52	0.43
Security	0	1	3	0	0	0	1	0	0	0.71	0.39
Performance	15	48	3	33	12	6	41	8	10	0.23	0.82
Stability	5	6	1	4	57	0	3	2	1	0.75	0.76
Feature	4	3	0	1	1	37	3	1	1	0.75	0.79
None	2	8	1	0	2	1	39	0	2	0.70	0.42
Mult.(corr.)	11	18	2	3	6	11	20	5			
Mult.(false)	4	7	2	0	2	2	6		1		
Macro avg. Accuracy										0.61	0.58
										0.55	

**Fig. 4.** Distribution of the judgments assigned in P3→P4 (combined score) and in P3' compared to the distribution of classes in the gold standard. Individual crowd workers could not assign multiple tags. See the online appendix (Groen et al., 2025) for details.



(a) Distribution of tags in P1.



(b) Distribution of tags in P2.

Fig. 5. Distribution of *helpful* vs. *useless* tags in Phases P1 and P2. The horizontal line indicates the share of *useless* tags in the gold standard of each phase.

Table 10

Accuracy of the crowd workers by agreement per class for Phases P1 & P2.

Correctness	P1 (N, %)		P2 (N, %)	
3 of 3	684	70.01	327	53.87
2 of 3	175	17.91	171	28.17
1 of 3	77	7.88	74	12.19
0 of 3	41	4.20	35	5.77
Total	977		607	
Average (%)	84.70		76.72	

3.3.1. General statistics

Overall, we used fairly long prompts for this study ($M = 12,393$ characters, $SD = 3556$), which included the data to classify. The processing time per batch averaged 1:08 min ($SD = 36$ s). ChatGPT required an average of 1.76 s per decision, and 0.71 decisions per second. We provide detailed analyses and all data in our online appendix (Groen et al., 2025).

3.3.2. LLM classifications for P1 & P2

Compared to the gold standards for P1 & P2, there is an apparent class imbalance toward classifying items as *helpful*, with ChatGPT 4 Legacy, in particular, predicting *useless* far less often, as shown in Fig. 5. While the LLM-based classifiers achieved high recall but limited precision on *helpful* in both phases, precision for *useless* was high, but recall was limited. We found strong main effects and several significant interaction effects in P1 for our study's three experimental factors:

1. *LLM* (4 vs. 4o): We determined a significantly better performance for ChatGPT 4o than for ChatGPT 4 Legacy ($\beta = 1.04$, $SE = 0.10$, $z = 10.07$, $p < .001$, $OR = 2.83$), which suggests a clear improvement of the updated model over its predecessor, making it 2.8× more likely to achieve better results on this task.
2. *Prompt type* (Eng vs. Kyō): The Kyōryoku prompt outperformed the engineered prompt ($\beta = 0.52$, $SE = 0.10$, $z = 5.40$, $p < .001$, $OR = 1.67$), indicating that the structure of the prompt affects performance, with the prompt that was not specifically engineered for ChatGPT being 1.7× more likely to achieve better results.
3. *Learning strategy* (Few vs. Zer): The role of a few-shot vs. zero-shot prompting strategy had no significant effect ($\beta = -0.10$, $SE = 0.09$, $z = -1.14$, $p = .26$, $OR = 0.90$), with the odds of omitting examples to negatively affect performance being just 10% in P1.

In addition, we found that the impact of the learning strategy depended on the LLM used, as evidenced by a significant interaction

effect for prompt type \times learning strategy ($\beta = -0.29$, $SE = 0.13$, $z = -2.22$, $p = .03$, $OR = 0.75$), and a three-way effect ($\beta = 0.39$, $SE = 0.21$, $z = 1.85$, $p = .06$, $OR = 1.48$). Further analysis suggested that ChatGPT 4 Legacy possibly benefited more from the Kyōryoku prompt, while ChatGPT 4o appears to have benefited from the examples of the few-shot learning strategy for both prompt types.

Table 12 shows the confusion matrix for the eight conditions in P1. ChatGPT 4 Legacy achieved better precision, while ChatGPT 4o achieved better recall. The $\langle Kyō, Few, 4o \rangle$ condition scored best with the highest recall score. The $\langle Eng, Zer, 4 \rangle$ condition scored poorest due to limited recall. Fig. 6(a) shows that both LLMs returned few false positives (FPs), with an area under the ROC curve (ROC-AUC) of 0.93 suggesting very good classifier performance in P1. However, the experimental factors LLM and prompt type appear to have had an effect on classifier performance. ChatGPT 4o consistently sacrificed specificity ($1 - FP$ rate) for sensitivity (TP rate), making it the best choice if a high TP rate is preferred (e.g., above 60%), while ChatGPT 4 Legacy achieved higher specificity at lower sensitivity, and would be the better choice for instances where the FP rate should be low (e.g., below 10%). Regarding prompt type, within the performance of each LLM, the engineered prompts always achieved higher specificity but equal or lower sensitivity compared to the Kyōryoku prompt.

In P2, we observed a much less clear pattern, with the bias toward tagging items as *helpful* being most pronounced in the $\langle Eng, Few \rangle$ and $\langle Eng, Zer \rangle$ conditions. ChatGPT 4o continued to score better than ChatGPT 4 Legacy ($\beta = 0.35$, $SE = 0.08$, $z = 4.24$, $p < .001$, $OR = 1.42$), but there was only a trend toward significance for performance using the Kyōryoku prompt versus the engineered prompt ($\beta = 0.15$, $SE = 0.08$, $z = 1.82$, $p = .07$, $OR = 1.16$). Neither the learning strategy nor any of the interactions were significant.

Table 13 shows the confusion matrix for the eight conditions in P2. The performance of both LLMs decreased, but while ChatGPT 4 Legacy achieved slightly lower precision and similar recall compared to P1, the performance of ChatGPT 4o saw a stark decrease in both precision and recall. The $\langle Kyō, Few, 4o \rangle$ and $\langle Kyō, Zer, 4o \rangle$ conditions scored best, both with good precision and fair recall. The poorest-performing condition was $\langle Eng, Zer, 4 \rangle$ due to limited recall, as in P1. Because P2 contained sentences from online user reviews of which at least one sentence was found *helpful* in P1, the task could be considered somewhat more difficult. The ROC plot in Fig. 6(b) reflects the decrease in performance in P2, with a lower ROC-AUC of 0.90. The difference in performance of the two LLMs and the conditions for each LLM are less distinct than in P1. However, there are similar patterns of a relatively low FP rate overall, higher sensitivity but lower specificity for ChatGPT 4o compared to ChatGPT 4 Legacy, and higher sensitivity but mostly lower specificity for the Kyōryoku prompt compared to the engineered prompt.

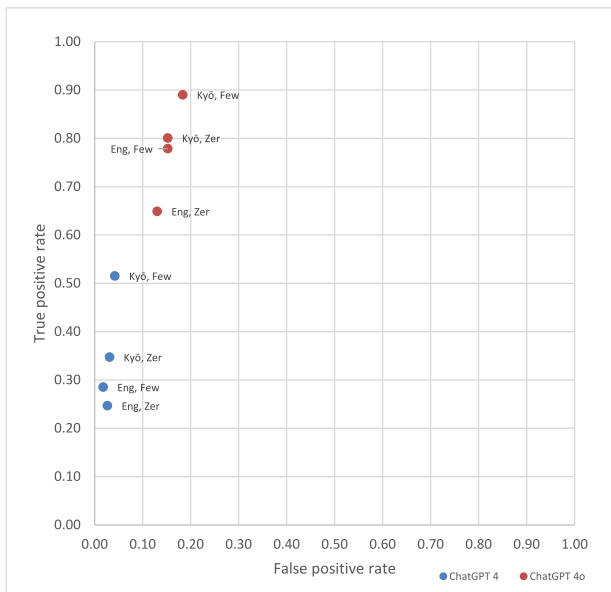
In Phases P1 & P2, we see that ChatGPT 4 Legacy consistently favored classifying items as *helpful* over *useless*. This might suggest

Table 11
Accuracy of the crowd workers by agreement per class for Phases P3 & P4.

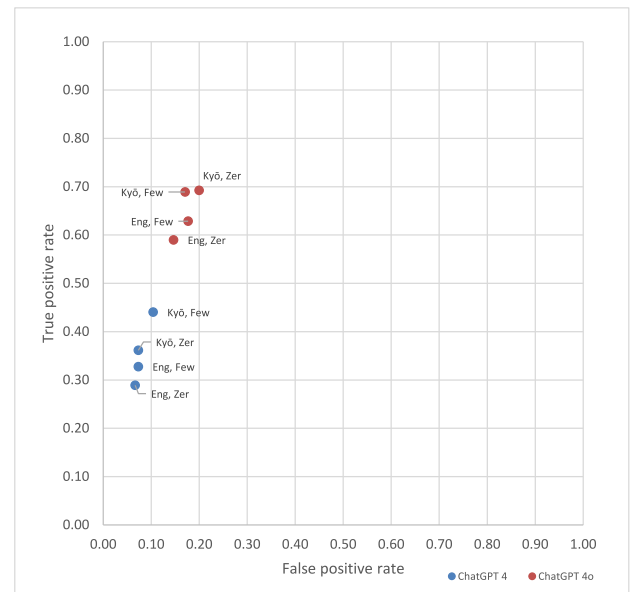
Phase	Agreement	Frequency (N, %)		Correct	Incorrect	Accuracy
P3	Six out of six	72	12.61	61	11	0.85
	Five out of six	117	20.49	89	28	0.76
	Four out of six	145	25.39	91	54	0.63
	Three out of six	157	27.50	70	87	0.45
	Two out of six	80	14.01	47	33	0.59
	Total	571	100.00	358	213	0.63
P4	Six out of six	28	14	26	2	0.93
	Five out of six	43	21.5	34	9	0.79
	Four out of six	57	28.5	35	22	0.61
	Three out of six	68	34	44	24	0.65
	Two out of six	4	2	4	0	1.00
	Total	200	100.00	143	57	0.72
P3'	Six out of six	53	9.28	48	5	0.91
	Five out of six	64	11.21	42	22	0.66
	Four out of six	147	25.74	81	66	0.55
	Three out of six	205	35.9	77	128	0.38
	Two out of six	101	17.69	64	37	0.63
	No agreement	1	0.18	0	1	0.00
	Total	571	100.00	312	259	0.55

Table 12
Confusion matrix comparing the results of the LLMs to the gold standard for P1. The number of positives and negatives is based on the gold standard. The reported values indicate the ability of LLMs to correctly identify items as *useless*.

Condition	Positives: <i>useless</i> (gold std.)	Negatives: <i>helpful</i> (gold std.)	True positives	True negatives	False positives	False negatives	Precision	Recall
Eng, Few, 4	516	484	152	472	12	364	0.93	0.29
Eng, Zer, 4	516	484	135	472	12	381	0.92	0.26
Eng, Few, 4o	516	484	415	404	80	101	0.84	0.80
Eng, Zer, 4o	516	484	347	417	67	169	0.84	0.67
Kyō, Few, 4	516	484	277	460	24	239	0.92	0.54
Kyō, Zer, 4	516	484	185	465	19	331	0.91	0.36
Kyō, Few, 4o	516	484	455	369	115	61	0.80	0.88
Kyō, Zer, 4o	516	484	415	392	92	101	0.82	0.80
Total	4 128	3 872	2 381	3 451	421	1 747	0.85	0.58



(a) ROC plot for P1.



(b) ROC plot for P2.

Fig. 6. Receiver Operating Characteristic (ROC) plots for Phases P1 and P2.

Table 13

Confusion matrix comparing the results of the LLMs to the gold standard for P2. The number of positives and negatives is based on the gold standard. The reported values indicate the ability of LLMs to correctly identify items as *useless*.

Condition	Positives: <i>useless</i> (gold std.)	Negatives: <i>helpful</i> (gold std.)	True positives	True negatives	False positives	False negatives	Precision	Recall
Eng, Few, 4	595	752	195	697	55	400	0.78	0.33
Eng, Zer, 4	595	752	172	702	50	423	0.77	0.20
Eng, Few, 4o	595	752	374	619	133	221	0.74	0.63
Eng, Zer, 4o	595	752	351	642	110	244	0.76	0.50
Kyō, Few, 4	595	752	262	674	78	333	0.77	0.44
Kyō, Zer, 4	595	752	215	697	55	380	0.80	0.36
Kyō, Few, 4o	595	752	410	624	128	185	0.76	0.69
Kyō, Zer, 4o	595	752	412	602	150	183	0.73	0.69
Total	4760	6016	2391	5257	759	2369	0.76	0.50

Table 14

Distribution of the judgments assigned by the eight LLM conditions in P3'. Note that all multiple tags are counted separately, causing the sum of judgments for each condition and the gold standard to exceed 624, the number of sentences in the gold standard.

Tag	Gold	Eng, Few, 4	Eng, Zer, 4	Eng, Few, 4o	Eng, Zer, 4o	Kyō, Few, 4	Kyō, Zer, 4	Kyō, Few, 4o	Kyō, Zer, 4o	Average
Compatibility	82	55	42	97	83	33	32	55	111	64
User-friendliness	153	50	57	85	116	40	75	101	88	77
Security	28	10	15	11	14	10	13	12	17	13
Performance	62	69	59	80	78	49	82	97	78	74
Stability	117	142	158	144	150	162	164	151	132	50
Feature request	79	99	91	125	126	105	96	120	120	110
None	147	218	212	108	86	233	169	100	94	153
Total	668	643	634	650	653	632	631	636	640	640

that ChatGPT 4 Legacy has a tendency to prioritize precision over recall, which resulted in high precision but limited recall for the *useless* classification. This effect was even more pronounced in the conditions that used the engineered prompt. In P1, ChatGPT 4 Legacy even achieved near-perfect recall on *helpful*, but had only limited precision. ChatGPT 4o, on the other hand, was able to achieve a more distributed assignment of tags especially in P1, where in the <Kyō, Few> and <Kyō, Zer> conditions, it even classified more items as *useless* than *helpful*, achieving a score of about 0.80 on both precision and recall in P1. However, its performance deteriorated in P2, achieving a score of roughly 0.70 on average, and displaying a pattern more similar to that of ChatGPT 4 Legacy in terms of favoring the classification of items as *helpful* over *useless*.

3.3.3. LLM Classifications for P3'

Table 14 shows the distribution of judgments for the eight LLM conditions in Phase P3'. Compared to the gold standard, *user-friendliness* and *security* received far fewer classifications, and *feature request* and *stability* far more classifications. The classification behavior of ChatGPT 4 Legacy and ChatGPT 4o for the other classes varied greatly. ChatGPT 4 Legacy classified far more and ChatGPT 4o somewhat fewer items as *none* compared to the gold standard, and ChatGPT 4 Legacy classified *compatibility* less often, while ChatGPT 4o classified *performance* more often.

Both ChatGPT 4 Legacy and ChatGPT 4o achieved considerably lower scores in P3' compared to P1 and P2, with a significant interaction effect for phase \times LLM ($\beta = -0.38, SE = 0.04, z = -9.60, p < .001, OR = 0.69$), as well as main effects for phase ($\beta = -0.06, SE = 0.03, z = -2.38, p = .02, OR = 0.94$) and LLM ($\beta = 1.23, SE = 0.08, z = 15.15, p < .001, OR = 3.43$). The effect sizes suggest that although ChatGPT 4o was 3.4 \times more likely to perform better than ChatGPT 4 Legacy, the effect of the LLMs on performance decreased by 31% in each consecutive phase, with the odds of the outcome decreasing by 6% per phase; a number that is only this low because ChatGPT 4 Legacy performed better in P2 than in P1. Fig. 7 illustrates this effect, with the curve for ChatGPT 4o showing a clear decrease in performance.

That none of the main and interaction effects in P3' were significant suggests that ChatGPT 4o did not have an advantage over ChatGPT 4 Legacy in this classification task. Although the main effect for the

LLMs did show a trend toward significance, with 19% odds of improved performance ($\beta = 0.17, SE = 0.12, z = 1.46, p = .14, OR = 1.19$), the accuracy of the eight conditions hardly differed, ranging from 0.58 to 0.64, with both prompt type and learning strategy having negligible effects. Prediction by majority vote achieved the highest accuracy score (0.67), although its advantage over the best-performing condition was not statistically significant ($\beta = -0.14, SE = 0.12, z = -1.13, p = .26, OR = 0.87$).

Table 15 summarizes the performance of the eight LLM conditions. ChatGPT 4 Legacy achieved low precision on *none*, among other things because it often assigns this tag to items that should be classified as *user-friendliness*. ChatGPT 4o, on the other hand, often incorrectly classified sentences as *compatibility*, independent of the frequency with which it assigned this tag. Prediction based on majority vote approached or even equaled the scores of the best-performing conditions while compensating for the poor predictions of the least well-performing conditions.

In individual classes, the classifiers were best at predicting *stability*, with fair precision and good recall, and *feature request*, with fair precision and good recall, caused by overclassification of this class. Recall was the greatest detrimental factor for the worst-performing tags—*compatibility*, *user-friendliness*, and *security*—due to underclassification (with the exception of ChatGPT 4o on *compatibility*). *Compatibility* was frequently misclassified as *none* by ChatGPT 4 Legacy and as *user-friendliness* or *feature request* by ChatGPT 4o; *user-friendliness* was misclassified into several different tags, and *security* was often misclassified as *stability*. Given a similar observation regarding *user-friendliness* with RQ2, this could suggest that the class may be more difficult to distinguish from other classes, the definition we provided is insufficient, or the LLMs construed a notion of *user-friendliness* that deviates from the definition and examples we provided. The majority vote prediction on individual classes was equal or superior to the best-performing condition. While Mizrahi et al. (2024) recommend averaging the output using multiple prompts, which inevitably causes the highest scores to be averaged out by the lowest scores, our evidence suggests that a majority vote has the potential to compensate for the poor precision or recall of individual conditions.

Table 16 shows the accuracy of the LLMs based on similarity between judgments. Note that here, we are comparing eight conditions with three different experimental factors: prompt type, learning

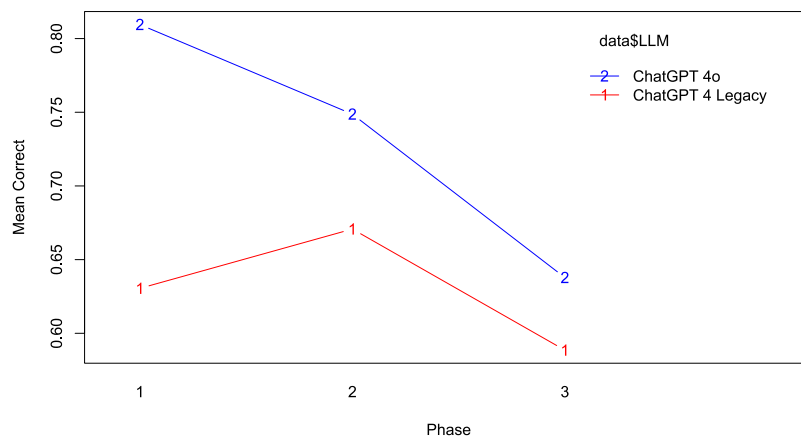


Fig. 7. Interaction effect for LLM across Phases P1, P2, and P3', showing a clear decrease in performance by ChatGPT 4o, and a less clear pattern but obviously reduced performance in Phase P3' by ChatGPT 4 Legacy.

Table 15

Performance per class by the eight LLM conditions in P3'. The full confusion matrices were omitted for space reasons and are provided in the online appendix (Groen et al., 2025). Red and green indicate the highest and lowest precision and recall value by class, i.e., compared horizontally.

	Eng, Few, 4		Eng, Zer, 4		Eng, Few, 4o		Eng, Zer, 4o		Kyō, Few, 4		Kyō, Zer, 4		Kyō, Few, 4o		Kyō, Zer, 4o		Majority vote	
	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R
Compatibility	0.52	0.37	0.61	0.32	0.48	0.56	0.47	0.46	0.68	0.27	0.80	0.31	0.54	0.37	0.46	0.61	0.64	0.45
User-friendliness	0.83	0.28	0.70	0.29	0.78	0.45	0.62	0.50	0.80	0.22	0.60	0.31	0.69	0.46	0.71	0.43	0.82	0.45
Security	0.60	0.33	0.64	0.41	0.64	0.33	0.67	0.41	0.70	0.35	0.69	0.35	0.67	0.29	0.69	0.35	0.82	0.47
Performance	0.53	0.65	0.60	0.58	0.58	0.80	0.60	0.85	0.62	0.53	0.49	0.73	0.48	0.86	0.60	0.87	0.59	0.83
Stability	0.69	0.84	0.65	0.89	0.72	0.88	0.71	0.91	0.62	0.88	0.62	0.87	0.68	0.88	0.74	0.86	0.66	0.91
Feature request	0.66	0.93	0.69	0.87	0.55	0.91	0.57	0.96	0.64	0.91	0.66	0.86	0.56	0.92	0.57	0.94	0.68	0.94
None	0.49	0.73	0.47	0.67	0.72	0.53	0.78	0.46	0.45	0.71	0.54	0.62	0.78	0.53	0.77	0.49	0.66	0.67
Macro average	0.62	0.59	0.62	0.58	0.64	0.64	0.63	0.65	0.64	0.55	0.63	0.58	0.63	0.62	0.65	0.65	0.69	0.67
Accuracy	0.60		0.59		0.64		0.64		0.58		0.59		0.63		0.64		0.67	

Table 16

Accuracy of the LLM conditions by agreement on one class for P3'.

Agreement	Freq (N, %)	Correct	Incorrect	Accuracy
Eight out of eight	269 43.11	230	39	0.86
Seven out of eight	63 10.10	46	17	0.73
Six out of eight	72 11.54	36	36	0.50
Five out of eight	109 17.47	53	56	0.49
Four out of eight	85 13.62	42	43	0.49
Three out of eight	25 4.01	12	13	0.48
Two out of eight	1 0.16	1	0	1.00
Total	624 100.00	420	204	0.67

strategy, and LLM. Agreement between multiple instances of the same condition would likely have been higher. The large share of perfect and near-perfect agreement suggests strong coherence between instances of ChatGPT despite different configurations. When all eight classifiers were in full agreement, they were accurate in 85.5% of the cases. Interestingly, accuracy quickly decreased, with three to six out of eight judgments in agreement all being at or just below 50% accuracy. The one outlier was the sentence “If you upgrade your power-ups to level 8, you will not be able to add them to your current loadout”. It had several different judgments given by a maximum of two classifiers, of which one answer was correct, causing the judgment to be counted as correct.

3.4. Results for RQ4 — Comparison

The crowd workers and ChatGPT Best (the best-performing ChatGPT condition) in P1 <Kyō, Few, 4o> performed similarly well ($\beta = -0.21, SE = 0.13, z = -1.59, p = .25, OR = 0.81$) and significantly better than ChatGPT Majority Vote, with the crowdsourced prediction being 2.1× more likely to perform better ($\beta = -0.75, SE = 0.13, z = -5.97, p < .001, OR = 0.47$) and ChatGPT Best 1.7× ($\beta = -0.54, SE =$

0.12, $z = -4.50$, $p < .001$, $OR = 0.58$). The precision and recall scores for the three classifier conditions for P1 in Table 17 show that ChatGPT Majority Vote outperformed the crowd on precision for *helpful* and recall for *useless*, but that its other scores did not match those of the other conditions.

In P2, the performance of the two best ChatGPT conditions was very similar, but our analysis showed that $\langle \text{Ky}\ddot{o}, \text{Few}, 4o \rangle$ had a slightly higher mean score than $\langle \text{Ky}\ddot{o}, \text{Zer}, 4o \rangle$. In this phase, the crowd achieved significantly higher scores compared to ChatGPT, being $1.8\times$ more likely to perform better than ChatGPT Best ($\beta = -0.60, SE = 0.10, z = -5.91, p < .001, OR = 0.55$), and $2.2\times$ better than ChatGPT Majority Vote ($\beta = -0.78, SE = 0.10, z = -7.75, p < .001, OR = 0.46$). The precision and recall scores for the three classifier conditions for P2 in Table 17 show that the crowd outperformed ChatGPT overall. Despite its poor precision on *useless*, ChatGPT Majority Vote was not significantly worse than ChatGPT Best, although we did observe a trend toward significance ($\beta = -0.18, SE = 0.09, z = -2.03, p = .11, OR = 0.83$), particularly due to the former’s poor precision on *useless*.

Table 18 shows the precision and recall scores for the five classifier pipelines considered for our comparison in P3 & P4. We excluded the *feature request* class because it was not part of the LP experiment of

Table 17

Performance per class for the crowdsourced tasks and LLM conditions compared to the gold standards for P1 & P2. Green indicates the highest precision and recall value by class and phase, i.e., compared horizontally for P1 & P2.

	Crowd		Phase P1		ChatGPT maj.		Crowd		Phase P2		ChatGPT maj.	
	P	R	P	R	P	R	P	R	P	R	P	R
Helpful	0.93	0.83	0.82	0.86	0.96	0.70	0.91	0.85	0.83	0.77	0.92	0.71
Useless	0.84	0.93	0.89	0.85	0.60	0.94	0.80	0.87	0.69	0.76	0.48	0.81
Macro average	0.89	0.88	0.85	0.86	0.78	0.82	0.85	0.86	0.76	0.77	0.70	0.76
Accuracy	0.88		0.86		0.78		0.86		0.77		0.73	

Table 18

Performance per class by the five selected conditions for RQ4 compared to the gold standard for P3 & P4. Red and green indicate the highest and lowest precision and recall value by class, i.e., compared horizontally.

	Language patterns		Crowd P3→P4		Crowd P3'		ChatGPT best		ChatGPT maj. vote	
	P	R	P	R	P	R	P	R	P	R
Compatibility	0.40	0.06	0.54	0.87	0.60	0.47	0.47	0.67	0.68	0.48
User-friendliness	0.81	0.10	0.50	0.61	0.54	0.44	0.71	0.52	0.82	0.49
Security	0.42	0.24	0.20	1.00	0.71	0.39	0.69	0.43	0.82	0.53
Performance	0.47	0.13	0.57	0.49	0.24	0.84	0.60	0.87	0.59	0.83
Stability	0.97	0.50	0.90	0.64	0.75	0.77	0.74	0.88	0.66	0.92
None	0.31	0.95	0.70	0.62	0.71	0.40	0.77	0.53	0.66	0.71
Macro average	0.56	0.33	0.57	0.71	0.59	0.55	0.66	0.65	0.71	0.66
Accuracy	0.41		0.64		0.53		0.66		0.68	

RQ1 and does not constitute a quality characteristic. Although this did slightly change the precision and recall values compared to the data reported above, it allows for a fairer comparison. With accuracy ranging from 0.41 to 0.68, the effect of the classifier was obviously significant ($\beta = 0.22, SE = 0.03, z = 8.08, p < .001, OR = 1.24$). Unsurprisingly, Tukey's post-hoc test showed that the classification based on the LP-based approach was significantly worse ($p < 0.001$) than all other classifiers. With RQ2, we had already established that Crowd P3→P4 had achieved statistically significantly better results than Crowd P3', but Crowd P3' was also outperformed by ChatGPT Best ($\beta = 0.40, SE = 0.12, z = 3.36, p = .007, OR = 1.49$) and ChatGPT Majority Vote ($\beta = 0.53, SE = 0.12, z = 4.45, p < .001, OR = 1.70$). With RQ3, we determined that ChatGPT Majority Vote was not statistically significantly different from ChatGPT Best. In this comparison, we also observed no statistically significant difference of both these results to those of the crowd in Phases P3→P4. This shows that in terms of overall performance, the three best classifiers did not differ from each other.

Table 19 summarizes the results of our per-class comparisons for P3 & P4. There were significant differences between the conditions for *user-friendliness*, *performance*, *stability*, and *none*, and a modest significance for *compatibility*. The findings per class are as follows:

- Pairwise comparisons for *user-friendliness* showed that ChatGPT Majority Vote differed significantly from all classifiers except ChatGPT Best, and that the poorest-performing classifier, Crowd P3', scored significantly lower than both ChatGPT classifiers.
- For *performance*, the difference between classifiers was the most pronounced of all classes because Crowd P3' scored considerably lower than all other classifiers; ChatGPT Best was better than both crowd classifiers.
- For *stability*, Crowd P3→P4 had significant pairwise differences compared to all other classifiers. For *none*, the difference between classifiers was very pronounced because the worst-performing classifier, Language Patterns, scored significantly lower than all other classifiers. ChatGPT Best differed significantly from Language Patterns and Crowd P3→P4.
- For *compatibility*, none of the pairwise comparisons were significant, and the prediction by Crowd P3→P4 did not differ significantly from the mean.

- *Security* was the only class for which there were no significant differences, and the prediction by ChatGPT Majority Vote was also not significantly different from that of the other classifiers.

4. Discussion of findings

This section discusses the implications of our findings according to the four research questions.

RQ1. How to adapt the NFR Method in order to derive LPs for the identification of quality characteristics in online user feedback?

To answer RQ1, we determined the feasibility of adapting the NFR Method (Dörr, 2011) and, in a sub-question, the accuracy of the outcomes of applying LPs. Regarding the adaptation, we obtained mixed results. Expert workshops were a suitable means of eliciting a substantial number of K&Ps related to quality characteristics and subcharacteristics—a level of granularity that research works usually do not address (Groen et al., 2017a). We also managed to encode LPs from the K&Ps. Often, several K&Ps were combined into one LP, but because negations and antonyms were placed into different LPs, about as many LPs were created as the number of K&P pairs that were obtained. However, LPs have been found to be inefficient because the manual vetting process involves substantial human labor, which outweighs any benefits of reusability and potentially high precision. In addition, they achieved low recall, partially because the coverage of relevant content based on the K&Ps elicited from just one workshop is limited.

RQ1a was aimed at assessing the effectiveness of the LPs based on quality-related keywords in identifying quality characteristics. The initial set of LPs achieved an average precision of 65%, which increased to 87% after refinements were made based on the seen data. Our validation over an independent test set showed a drastic reduction in precision to 56% on average, with two classes retaining high precision. Despite the fact that the LPs maintained good precision on *usability*, they achieved the second-lowest recall score in this class. However, *reliability* achieved near-perfect precision with good recall. This is an important finding because reliability issues often impact larger groups

Table 19

Statistics for the performance of the five selected conditions for RQ4 per class for P3'. The Friedman Test (χ^2 ; $df = 4$, $N = 511$) indicates whether or not there is a significant difference between classifiers; the binomial GLM (β) measures whether the highest mean reported in the table is significantly different from the average for all classifiers.

Class	χ^2	p	Highest mean	β	SE	z	p	OR
Compatibility	13.0	<0.01	Crowd P3→P4	0.26	0.15	1.80	0.07	1.30
User-friendliness	29.5	<0.001	ChatGPT maj.	0.40	0.12	3.40	<0.001	1.49
Security	4.2	0.374	ChatGPT maj.	0.23	0.24	0.95	0.34	1.26
Performance	265.0	<0.001	ChatGPT best	0.97	0.16	5.95	<0.001	2.63
Stability	28.5	<0.001	ChatGPT best	0.38	0.15	2.48	<0.01	1.47
None	253.0	<0.001	ChatGPT best	0.81	0.12	6.97	<0.001	2.25

of end-users; sentences pertaining to this class mostly report on the app crashing, freezing, hanging, breaking, or quitting. Tested against our gold standard, the LP-based approach typically achieved low recall, which suggests that LPs are very selective in their prediction of quality characteristics. This is due to the heterogeneity of online user feedback, where laypeople usually do not describe the circumstances in sufficient detail. Moreover, an approach that rigidly matches a combination of words is unable to infer cues from the context that could help resolve ambiguities, and consequently will have difficulty correctly determining the affected quality. For example, the phrase “app does not work” could imply a problem with *installability*, *reliability*, *compatibility*, *interoperability*, or *security*, or maybe reflect the user failing to understand how to use the app. This may also be due to end-users perceiving the software exclusively from a front-end perspective.

Answer to RQ1: Expert workshops are suitable for producing a keyword meta-model for software product quality characteristics, but the encoding of language patterns is a time-intensive task. The performance of the LPs when classifying online user feedback was poor, and the LPs were only able to classify sentences about *reliability* fairly well.

RQ2. How to design micro-tasks so that lay workers recruited through crowdsourcing can achieve good results classifying quality aspects in online user feedback?

For RQ2, we extended the research on Kyōryoku (van Vliet et al., 2020), a crowdsourcing method for eliciting user requirements from online user feedback. We investigated how the decomposition of a classification task into smaller, consecutive classification tasks affects the performance of crowd workers (RQ2a). In a single-group experiment, we attracted a crowd of 711 non-expert human annotators, 555 of whom eventually contributed to the classification of a random sample of 1000 app store reviews into requirements-relevant dimensions and quality aspects. We compared two conditions: one in which we presented all quality aspects at once (Phase P3'), and one in which we split the classification process into two micro-tasks with fewer quality aspects (Phases P3→P4).

We found that the crowd workers performed reasonably well on this complex task in the decomposed P3→P4 sequence, achieving an accuracy of 63% and 72%. The crowd workers in P3' performed considerably worse, with an accuracy of 55%. Decomposing the classification problem into several simpler tasks seems to lead to better crowd performance. However, there is also a downside to performing multiple tasks in sequence: Misclassified sentences (i.e., reduced precision) in a class are prevented from being correctly classified in a later phase, at the detriment of recall for the phases combined. In this study, the sentences that the crowd did not classify as *quality* in P3 (see Table 7) would not proceed to P4. In P3→P4, the severity of this problem appears to have been mitigated by a lower tendency to classify sentences as *none*, allowing more sentences to be classified in P4.

The observation that the larger number of tagging options in P3' negatively impacted the crowd's accuracy is in line with general research on manual annotation (Bayerl and Paul, 2011). In particular, the broad spectrum of choices added confusion and resulted in crowd

members defaulting to one tag when in doubt—specifically for *performance*—which reduced precision. Further analysis suggests that the crowd workers often made interpretations beyond the definitions given in the instructions. It appears that these interpretations particularly came into play for ambiguous sentences, i.e., in the case of multiple statements, contradictions, or unclear language. We also found evidence of this in the crowd workers' responses to our test questions, some of whom defended their dissenting judgment. This does not only occur among laypeople; while devising the gold standards, the researchers also observed how personal heuristics interfered with their classifications.

Crowd workers were generally able to differentiate between software product qualities represented by their classes (RQ2b), even if some occasionally got confused. The frequency and accuracy of each assigned class differed between phases, usually with different classes getting confused. Instances with perfect agreement corresponded to high accuracy scores. Although it is likely to assume that the highest agreement is achieved on the least ambiguous sentences, a configuration requiring a minimum number of crowd workers agreeing with one another might lead to more accurate results. Like other studies (e.g., Schenk and Guittard, 2011; Gilardi et al., 2023), we found that the performance of crowd workers deteriorated as task complexity increased. It is nevertheless possible to crowdsourcing complex tasks that could even be challenging to expert judges, with the best results obtained if they are decomposed into smaller, less complex ones with clear instructions.

Answer to RQ2: A crowd of lay workers is capable of annotating online user feedback according to requirements relevance and even functional and quality aspects. Clear instructions help to convey the required expert knowledge. The best performance we attained was through a series of micro-tasks, each focusing on only a few classes.

RQ3. How to design an LLM pipeline so that it can achieve good results in classifying quality aspects in online user feedback with little training data?

For RQ3, we performed a $2 \times 2 \times 2$ between-subjects experiment, with a pipeline that prompted two LLMs from the ChatGPT family, as well as the factors prompt type and learning strategy, resulting in eight conditions. Investigating how increasing task complexity affects the performance of LLMs (RQ3a), we found that ChatGPT achieved high precision but mostly limited recall for the binary classification task over online user reviews (Phase P1). For the somewhat more difficult binary classification task over individual sentences (P2), precision was still high, but recall worsened. The most difficult multiclass classification task over individual sentences (P3') resulted in mediocre precision and recall. The performance of the LLMs worsened as task complexity increased, resulting in a decrease in the average accuracy scores of the conditions across the three phases of our experiment. Which LLM was used had the greatest impact on performance, as ChatGPT 4o made significantly better predictions on simpler tasks than ChatGPT 4 Legacy. However, ChatGPT 4o's upper hand diminished across the three phases, while the performance of ChatGPT 4 Legacy only worsened gradually.

These findings suggest that LLMs are suitable but not perfect for this task.

Because some literature suggests that LLMs are likely to perform better with a specifically engineered prompt than with instructions not specifically engineered for LLMs (e.g., Brand, 2023; El-Hajjaji et al., 2024), we studied how the use of engineered prompts and examples affects the performance of LLMs (RQ3b). In Phases P1 and P2, the prompt type did seem to have an influence on performance, but surprisingly, it was a reverse effect. The Kyōryoku prompt not specifically engineered with and for ChatGPT tended to lead to better results than the specifically engineered one. This effect was significant for P1, and in P2 had a trend toward significance. In Phase P3', the prompt type neither positively nor negatively affected performance. Possible explanations for this finding include our limited experience with prompt engineering, despite the iterative approach taken to carefully align with ChatGPT, the shorter length of the Kyōryoku prompt allowing ChatGPT to find important guiding information more easily, or the notion that ChatGPT can be considered a layperson similar to a crowd worker, which was the intended target audience of the Kyōryoku instructions.

Regarding the role of learning strategy, the literature was inconclusive. In P1, there was an interaction effect for prompt type \times learning strategy, suggesting that the examples we provided mostly had no effect on performance, except helping ChatGPT 4 Legacy to perform slightly better in P1. Reynolds and McDonnell (2021) provide a good explanation of how to understand the role of few-shot examples. They argue that LLMs do not use prompts for learning, but instead for assigning the appropriate task location in the LLM. Consequently, the few-shot examples merely steer the retrieval closer to or further away from the correct knowledge space. In our research, the definitions for each class may have been sufficiently self-explanatory for the examples to affect task location.

Regarding performance by class, we observed a class imbalance for ChatGPT 4 Legacy on P1, and for both LLMs on P2, which mostly overclassified items as *helpful*. This pattern seems to resemble the phenomenon that ML-based classifiers applied to unseen data default to the positive answer when in doubt, prioritizing recall over precision. We are unsure why this happens when the mechanics are different; ML draws inferences from the data on which it was trained, while the decoder architecture of LLMs predicts the most likely response based on its model (Raffel et al., 2020). In our experiment, this entailed predicting which of the predefined class names should be returned. That ChatGPT 4o behaved differently in P1 than in P2 might be due to it erring on the side of caution; for example, in P2 after sentence splitting due to the lack of context from a full review. For the multiclass classification in P3', we found that both LLMs tended to overclassify *feature request* and *stability* and underclassify *user-friendliness* and *security*, resulting in low recall. The LLM conditions rarely achieved a balance between precision and recall for a single class. This suggests that in this complex task, both LLMs could not always correctly discern between the classes presented to them.

Answer to RQ3: LLMs are capable of performing classification tasks according to requirements relevance and functional and quality aspects. In our study, ChatGPT 4o mostly outperformed ChatGPT 4 Legacy. The choice of prompts sometimes improved performance, but providing examples (i.e., few-shot learning) had no effect.

RQ4. Which low-data approach is most suited for classifying quality aspects in online user feedback, considering effectiveness and trade-offs?

We compared the three approaches both quantitatively and qualitatively to identify the most suitable one for identifying quality aspects in online user feedback. The crowdsourced Kyōryoku method and the

LLM pipeline achieved comparable results on the classification problem of annotating quality aspects in online user feedback, and were much more accurate than the LP-based approach. This finding highlights that when appropriate measures are taken, this complex classification task can be performed by a hybrid system (crowdsourcing) or an expert system (LLMs) instead of being performed or supervised by domain experts.

Regarding the suitability for determining requirements relevance (RQ4a), the crowd and the best-performing ChatGPT condition performed equally well on P1, with the crowd achieving better recall and ChatGPT achieving better precision on *useless*. The crowd performed about equally well in Phase P2, but the best-performing ChatGPT condition performed worse, and ChatGPT's majority vote prediction was even less accurate, especially in terms of recall. This might suggest that human workers were better at considering sentences in isolation in P2, and that ChatGPT in part relied on the context provided in the complete online user reviews of P1, making it more difficult for it to draw correct inferences if this context was missing. Interestingly, the best ChatGPT condition in P1 and P2 was <Kyō, Few, 4o>, which happens to be the one prompted with the exact same instructions as the crowd workers received. Thus, ChatGPT 4o is capable of equating human workers on a classification task when given the same instructions, but the crowd was generally better at predicting whether or not online user feedback is requirements-relevant.

Our analysis of how suitable each approach is for distinguishing between the quality aspects (RQ4b) showed that classifying quality aspects in online user feedback remains a challenging task, given that it is a multiclass classification task based on an inherently arbitrary taxonomy of software qualities over intrinsically ambiguous online user feedback and performed without training data. The LP-based approach faced challenges due to its rigidity, but the approaches using crowdsourcing and LLMs also performed significantly worse on the classification task(s) in P3 & P4, although they could usually correctly predict the quality characteristics with at least good precision or recall. Three conditions were equally accurate: the crowd in the P3→P4 sequence, the best ChatGPT condition, and ChatGPT's majority vote prediction. The other two conditions in our comparison—the LP-based approach and the crowd in Phase P3'—did not perform well. The equal performance of ChatGPT and a crowd of human raters on this multiclass classification task suggests that the advantage of the knowledge-processing capacity and training of an LLM does not necessarily make it superior to a crowd of lay workers who have typical human cognitive limitations when predicting quality aspects in online user feedback.

Due to the different patterns we observed, we found no conclusive evidence that a specific classifier is the best on a per-class level. Our classifier conditions showed patterns that resembled those found with traditional ML classifiers. For example, as in our study, Kurtanović and Maalej (2017) obtained the highest precision and recall values for *performance* and the lowest for *operability* and *usability*, which in our study both belong to *user-friendliness*. Although the best-performing configurations for crowd workers and LLMs performed equally well, we did determine that each approach tends to behave differently. For example, Table 18 illustrates how precision and recall by Crowd P3→P4 was the inverse of ChatGPT Majority Vote (and, to a lesser extent, ChatGPT Best); while the former mostly excelled on recall and the latter mostly on precision, it was also the opposite for some other classes.

We see further improvement potential for all approaches. The crowd will likely perform better if more judgments are collected because greater agreement between crowd workers led to higher accuracy than agreement between LLM conditions. For the approach using LLM, a pipeline specifically designed for results based on a majority vote and the use of variations in prompts could further boost its accuracy. Even though the LP-based approach did not compare well to our other approaches, its ability to predict *stability* well and its potential for explainability and reusability might still be of use for settings that

prioritize predictable outcomes and transparent processes, which the other approaches cannot provide due to the personal heuristics employed by the crowd workers and the non-deterministic nature of LLMs. The LP-based approach could therefore benefit from automating the process of curating the lexicon of keywords and vetting the LPs using ML approaches or LLMs.

Answer to RQ4: Overall, crowdsourcing was better in determining requirements relevance. Crowdsourcing, the best-performing ChatGPT condition, and the majority vote prediction of ChatGPT classified quality aspects in online user feedback equally well, with no classifier clearly being the best in individual classes.

5. Threats to validity

Several threats to the validity of our studies' results should be considered. We will discuss the main validity threats according to the dimensions of construct, internal, external, and conclusion validity (cf. Wohlin et al., 2000).

Construct validity is concerned with the appropriateness of the research approach and the instrumentation to answer the research questions. For RQ1, we used an empirically validated method for eliciting quality requirements. The use of other linguistic approaches might have led to different outcomes. The evolution of the elicitation workshops for RQ1 due to the adjustments we made to improve the effectiveness of the instructions and templates could have changed the group dynamics, but these changes were minor. Although we had to switch from co-present to remote workshops due to the COVID-19 pandemic, we did not observe any impact on the workshops' dynamics or outcomes. We made the best possible effort to reduce bias and misinterpretation regarding quality characteristics, especially by dedicating time at the beginning of the workshop to making the participants' preexisting notions explicit and aligning their individual understandings of the quality characteristic with each other and with the ISO 25010 definition.

For RQ2, van Vliet et al. (2020) discusses the threats to validity in the design of Kyōryoku, with a particular focus on the lack of reference material and practical experience with crowdsourcing tasks. The finding that a sequence of micro-tasks leads to better results is based only on the P3→P4 sequence and would have been stronger if we had considered multiple sequences in which we varied the composition of classes.

The main discussion of construct validity concerns RQ3, because from an empirical point of view, our understanding of LLMs is still limited, and LLMs produce output that is non-deterministic. Several threats pertaining to study participants apply to measuring the behavior of LLMs as a generative entity, which in this work may include history, maturation, and testing effects. However, a threat of particular concern is *multiple-treatment interference*, which Campbell and Stanley (1963) describe as a threat through which generalizability is reduced when applying multiple treatments to the same respondents whenever the effects of prior treatments are not erasable. We actively mitigated the dependency of messages within a single ChatGPT conversation by performing all the prompt sessions separately, with ChatGPT's memory function disabled.

The choice to use ChatGPT in research is not undisputed, as it has both clear advantages and marked disadvantages. The purpose of this research was to investigate the best result that LLMs could potentially achieve at the time of the experiment, not an exhaustive comparative study between LLMs. Based on an extensive pretest (see the online appendix; Groen et al., 2025), we determined that ChatGPT was the most-researched and purportedly the most robust LLM for achieving that goal. From among the alternatives, it was the most capable in terms of handling our prompts and producing meaningful results. We encourage the evaluation of other LLMs through reproductions of this experiment. However, it is appropriate to elaborate on its disadvantages. We chose ChatGPT over prompting GPT through an API because

we had established that it had just received the feature of handling Excel files and that it was capable of processing larger amounts of text. In retrospect, using GPT might have been less time-consuming due to the amount of troubleshooting we needed to do. It could also have allowed certain hyper-parameter settings such as temperature or top-*k* to be altered, potentially affecting classification accuracy. We do not expect that the use of the Web interface affected our ability to program the experiment or avoid errors, but we did encounter response limits followed by cooldown periods and could not control the number of tokens, which is why we reported on the number of input characters instead. The use of ChatGPT is dependent on design decisions made by its developer OpenAI; its performance can be unpredictable due to its opaque evolution, which, in turn, also impairs its replication potential, but we also observed improvements being made to its processing speed. We estimated this limitation to be acceptable considering the response quality we established for the two state-of-the-art LLMs we compared. We consider that the two additional experimental factors of prompt type and learning strategy that we chose based on our pretest were appropriate for our comparison purposes.

Finally, LLMs are under scrutiny from an ethical stance, which includes concerns about regulations and data privacy, transparency, hallucinations, and harmfulness (Vogelsang, 2024). These should be taken into account in the context of other approaches considered. Although our dataset was curated from publicly available sources, we ensured that our data was anonymized (Groen and Ochs, 2019; van Vliet et al., 2020), and the fixed response options limited the potential for hallucination. Harmfulness particularly pertains to sustainability concerns about energy consumption. Here, we note that LLMs require a surge of energy, while our other approaches rely on power consumption over a prolonged period of time for vetting and deploying LPs or for performing the decentralized crowd work on hundreds of devices. Similarly, one should consider the trade-off between whether it is desirable to strain an expert with a cognitively demanding task, whether voluntary crowd workers from impoverished countries are exploited or if it allows them to earn additional income (cf. Sachs, 2005), and whether or not this outweighs sponsoring the shareholders of the big tech companies operating LLMs.

Internal validity addresses the probability that the same results will be obtained if a study is replicated. For RQ1, the variation between experts could cause workshops with different expert participants to produce different results. Nevertheless, this likely has only a limited effect on our interpretation of how well LPs can predict quality characteristics in online user feedback. Here, the primary weakness of using quality-related (combinations of) keywords recoded into LPs is the low recall that we were able to achieve. Because only predefined words were found, the LPs were inherently prone to missing relevant keywords that were not elicited. We had previously established that an approach in which keywords are extracted from an existing dataset can only partially mitigate this (see Study II in Groen et al., 2017a). For this work, the deviation from the NFR Method would have been a too large, and it would have interfered with our ability to draw inferences about expert contributions. In RQ2, we saw, in particular, that the crowd we had attracted for the tasks of P4 and P3' was different than that for P1, P2, and P3. Especially the fairness of remuneration was rated lower, even though we paid more. The crowd workers also found the tasks to be more difficult, even though P4 should have been easier than P3. This might suggest that crowd workers attracted at different points in time might have different results. The design of the approach in which each phase utilizes input from the preceding phases is likely to perpetuate errors further; especially regarding items not presented in P4 that the crowd workers in P3 failed to classify as *quality*. Despite this, the two consecutive Phases P3→P4 still achieved better results than Phase P3'. The quality of the results could possibly be improved by involving more crowd workers, although this also increases costs. The main threat to the internal validity of RQ3 is the non-deterministic

nature of LLMs, whose impact we discussed as part of construct validity. Although we found in our study that the performance of the LLMs was consistent across conditions—despite different experimental factors—their behavior might vary at different points in time.

An important determinant for the internal validity of all three studies is our dataset. Due to the rapid evolution of the app landscape, the results might have differed if more recent online user feedback had been used. However, we cannot gauge the impact of the age of the dataset because to our knowledge, there are no works that have replicated or updated the work by Pagano and Maalej (2013) that could give an indication of whether any characteristics of the content of online user feedback have shifted over time. Based on the analysis of the confusion matrices, we learned that more time needs to be spent on discussing the gold standards, both by reconciling differences in tagging and by verifying consistency. In van Vliet et al. (2020), this was partially mitigated by manually inspecting the differences and adding a lenient judgment. To avoid this unnecessary step, we refined the gold standards through iterative sessions among three authors for a more reliable comparison. Due to the choice of a bootstrapped random sample from a larger dataset, the different classes had varying numbers of sentences associated with them, resulting in the strength of the performance calculations to vary. This particularly affected *security*, which had the lowest number of sentences. We also performed extensive manual sanitation of the data upon finding contamination issues, in order to achieve a fair and reliable analysis. This resulted in the omission of approximately 15% of the data for RQ2, especially in P2, which could have negatively affected the performance values of the crowd.

Regarding RQ4, to ensure comparability between approaches, we had to make compromises, particularly in terms of their classification. The LP-based approach was limited to only one label. To keep the micro-tasks in Kyōryoku simple, each crowd member could provide just one label, from which a multi-label prediction arose in case of a tie in the majority vote. We instructed ChatGPT to assign a multi-label response only in case of a tie between possible labels, and ties in the majority vote prediction also resulted in multi-label predictions. We considered the multi-label items in the gold standard for P3 & P4 as a range of correct options, only one of which needed to be matched for a correct response. We leave explorations of multi-label classifications and attributing classes to smaller portions of the input text to further research. In addition, each approach had unique advantages and disadvantages that could affect performance, complicating the objective comparison of their strengths and challenges. The LP-based approach benefited from iterative refinements made to the LPs, whose high precision could not be offset at that stage against objective recall scores. In the Kyōryoku method, the crowd achieved good scores in Phase P4 due to the smaller number of classes, while ChatGPT benefited from having eight experimental conditions, as opposed to one for the LP-based approach and two in Kyōryoku (P3→P4 and P3'). We compensated for the weaknesses of each approach in different ways: In the Kyōryoku method, we compensated for the mental limitations of laypeople by breaking down the task into multiple classification stages; we provided the LLMs with more context and engineered prompts to prevent potential misinterpretations; and we iteratively refined and optimized the LPs.

External validity concerns the generalizability of the results. In RQ1, the workshops were conducted only once for each quality characteristic, but we argue that the discussion of the experts participating in the workshops sufficed to produce a well-composed keyword meta-model. The workshop participants were not frequent authors of online user reviews themselves, but a pretest had shown that computer laypeople provided fewer relevant keywords & phrases. Tizard et al. (2020) found that typical authors of user reviews have greater computer affinity, which provides some evidence that our participants can be considered sufficiently representative. In addition, our work shows similarities with other research in this direction, such as Cleland-Huang et al.

(2007b), who found that for several indicator words, it is the context that prescribes the subcharacteristic to which it should be attributed, and that the challenge accordingly lies in the boundaries between the quality characteristics. For example, do user reports about being able to do something with the app really “fast” relate to the user’s perception of the app’s *time behavior* or its *operability* in terms of the workflow?

In RQ2 and RQ3, the main threats to external validity are a logical consequence of the already described changes in the composition of the attracted crowd or the LLM used at a given point in time. Majority vote prediction might be more robust against changes in the performance of individual crowd workers or ChatGPT conditions. Based on the results of the pretest we performed for RQ3, we consider it unlikely that the pipeline we created will lead to the same results when prompting other LLMs than those tested.

Conclusion validity pertains to the credibility of the results obtained. We have done our best to ensure the quality of our data and interpret it using tabulations, visualizations, and appropriate statistical measures. Due to the newness of LLM technology, we have attempted to base our decisions on how to analyze it based on the available literature on evaluating classifiers, as part of which Dell’Anna et al. (2023) is recognized as an authoritative work (cf. Vogelsang, 2024). However, during this process, we faced important fundamental questions, such as whether the same LLM in a different condition should be analyzed using a repeated-measures test, or whether its non-deterministic nature is a reason not to do so. Due to the novelty of this domain, best practices for evaluating LLMs have not yet been developed, and these might eventually favor other metrics and tests than the ones we chose. We share our raw data in the online appendix (Groen et al., 2025), over which alternative statistical tests could be performed.

6. Related work

We will discuss the main related work by addressing crowdsourcing in RE (Section 6.1), automated analysis of online user feedback (Section 6.2), and LLMs in RE (Section 6.3).

6.1. Crowdsourcing in RE

The term *crowdsourcing* originated from the business sector, when Howe (2006) coined it as “the act of taking a job traditionally performed by a designated agent (usually an employee) and outsourcing it to an undefined, generally large group of people in the form of an open call”. In his comprehensive overview, Howe (2010) refers to this definition as “the White Paper Version”. In this traditional perception of crowdsourcing, a crowdsourcing platform offers a crowdsourcing task of a crowdsourcer to a crowd of workers (Hosseini et al., 2014a, 2015a). Our Kyōryoku method used for RQ2 belongs to the body of work that approaches crowdsourcing in this way. Expanding on his own definition, Howe (2010) provided an additional perception of crowdsourcing, where it is seen as “the application of Open Source principles to fields outside of software”. Howe refers to this pragmatic perception of crowdsourcing as “the Soundbyte Version”. No known work in RE specifically addresses the pragmatic perception, but it illustrates that there can be more than one notion of crowdsourcing. The majority of works in RE do, in fact, assume what could be described as a modernist perception of crowdsourcing, which considers crowdsourcing to involve any approach where information is obtained (i.e., “sourced”) from a crowd, typically not bound by Open Source principles and usually for profit. Neither Howe nor the literature provides an official definition for this view. Wang et al. (2019) provide an overview of approaches in RE that make use of “crowdsourced user feedback”, where crowdsourcing encompasses obtaining (i.e., “sourcing”) requirements from the crowd’s online user feedback. In RE, this is typically referred to using the umbrella term *Crowd-based Requirements Engineering* (CrowdRE), which encompasses all approaches to (semi-)automatically analyzing and classifying online user feedback for the specification and evolution

of software products in RE (Groen et al., 2015, 2017b; Glinz, 2019). Many of the approaches suggested in CrowdRE focus on the analysis of natural-language texts. Engaging the crowd mainly relies on a form of self-regulated social participation (Ali et al., 2012). We consider all three approaches proposed in this work to fall under the CrowdRE umbrella.

The early 2010s saw a drastic rise in the interest to adopt traditional crowdsourcing in many domains, including SE (for reviews, see Mao et al., 2017; Ghezzi et al., 2018). However, reports of poor results from crowdsourced activities might have caused the upswing to lose traction. Notably, Stol and Fitzgerald (2014) presented a case study in which software programming work was crowdsourced, where none of the promises of crowdsourcing—cost reduction, faster time-to-market, high quality, and creativity & open innovation—were attained (Fitzgerald, 2018). Their work shows that crowdsourcing might not be suited for tasks that are more complex or require coordination between developers.

Regarding crowdsourcing for RE, in their seminal work on this topic, Hosseini et al. (2014b) discuss five crowdsourcing approaches in RE, for which they suggested appropriate crowdsourcing platforms be developed: requirements-driven social adaptation, feedback-based RE (i.e., CrowdRE), stakeholders' discovery, requirements identification, and empirical validation. Moreover, they discuss ten crowdsourcing features and their associated quality attributes, which were expanded upon in Hosseini et al. (2015b) with a review of the challenges faced for each of these features. Khan et al. (2021) performed an interview study with 50 practitioners to elicit challenges when using crowdsourcing platforms in RE, finding challenges related to time, quality of work, culture, number of experts, communication and confidentiality, conformity between feedback and requirements, and degree of knowledge on crowdsourcing.

Among the tools proposed to use crowdsourcing in RE is *CrowdREquire* (Adepetu et al., 2012), an envisioned online platform to crowdsource requirements specifications in the form of a competition. Most other works that involve traditional crowdsourcing in RE tasks have been found to focus especially on mechanisms involving non-RE actors in requirements elicitation and validation (Groen, 2015; Hosseini et al., 2014a, 2015b). For example, *StakeSource* (Lim et al., 2010) is a research-based commercial crowdsourcing platform that supports stakeholder selection through peer recommendations and involves them in requirements elicitation and prioritization. It partly draws on the *Organizer & Promoter of Collaborative Ideas* (OPCI; Castro-Herrera et al., 2009), where stakeholders could co-create requirements in a forum that automatically clustered related ideas and used a recommender system to link stakeholders to potentially relevant threads. The *Crowd Requirement Rating Technique* (CrowdReRaT; Oluwatofunmi et al., 2020) was designed to crowdsource requirement ideation and rating, and a tool developed by Sari et al. (2021) proposed the use of crowd discussions in the elicitation process. More recent efforts have focused on identifying and selecting suitable stakeholders through crowdsourcing (Alamer and Alyahya, 2023; Delima et al., 2023).

Although the traditionalist perception of crowdsourcing does not include obtaining online user feedback indirectly, several approaches in this line of research have suggested ways to actively involve crowd workers in a crowdsourced process to solicit online user feedback, perform feedback analysis, or validate the results (Groen, 2015). Snijders et al. (2015) proposed *REfine*, a gamified CrowdRE approach that involves crowd members in gathering needs and in the subsequent requirements analysis and prioritization. This work was extended in the *Crowd-based Requirements Elicitation with User Stories* method (CREUS; Wouters et al., 2022), which gathers ideas from the crowd using the user story notation through four phases: preparation, idea generation, refinement, and response and execution. The initial ideas presented for *Requirements Bazaar* envisioned a similar solution for eliciting requirements through a platform, but additionally included considerations for requirements negotiation (Renzel et al., 2013). The *Crowd-Annotated*

Feedback Technique (CRAFT; Hosseini et al., 2017) crowdsources the classification of app store feedback according to RE-related dimensions. Crowd members act as taggers, who select a particular fragment of online user feedback, pick a category and a (sub-)classification, and provide an assessment of the feedback quality and their own confidence level. Although CRAFT was designed as an in-situ mechanism and did not offer remuneration as an incentive, its validation demonstrated the possibility of outsourcing online user feedback annotation to a crowd, which we later built on when designing *Kyōryoku* to crowdsource online user feedback classification through micro-tasks (van Vliet et al., 2020, see RQ2), which is the basis for RQ2. Aside from our work, the only other example of research using an RE-related crowdsourced annotation task is Stanik et al. (2019), who based part of their ML and DL training set on 10,000 English and 15,000 Italian tweets annotated by the crowd. Our work in RQ2 focuses on how this annotation task can be configured optimally.

6.2. Online user feedback analysis in RE

In RE, natural language processing (NLP) has mainly been concerned with identifying content in documents for RE tasks (Ferrari and Ginde, 2025). In our work, we applied a linguistic approach for RQ1 and performed requirements classification with all three approaches considered.

Linguistic approaches. Performing text mining using a linguistic approach involves identifying patterns in utterances that point to a possible use or intention (Jurafsky and Martin, 2009). Although most of the recent literature has favored the use of traditional ML (for a review, see Santos et al., 2019a), linguistic approaches such as keyword tracing have successfully identified requirements in natural-language texts. These are rooted in strategies that far precede the evolution of online user feedback analysis. Cysneiros et al. (2001) proposed a strategy to integrate non-functional requirements (NFRs) into conceptual models by constructing a Language Extended Lexicon consisting of a shared application vocabulary that links representations such as requirements and models. This technique could be used to elicit, analyze, and trace NFRs. This makes domain ontologies suitable for requirements management activities such as glossary development or trace link vetting (Dermeval et al., 2016). For example, the *OpenReq Dependency Detection* tool (OpenReq-DD; Motger et al., 2019) automatically generates ontologies and clusters keywords to detect requirements dependencies (Deshpande et al., 2020). For RQ1, we developed a keyword meta-model for software product quality characteristics in the form of a lexicon.

The first application of linguistic techniques to online user feedback analysis was the *Mobile App Repository Analyzer* (MARA; Iacob and Harrison, 2013), which identifies feature requests in text fragments of user reviews using language patterns and aggregates results. Regular expressions have proven useful for extracting positive feedback, negative feedback, and feature requests from user reviews (Fu et al., 2013; Iacob and Harrison, 2013). Panichella et al. (2015) derived 246 *recurring linguistic patterns* through a manual inspection of 500 user reviews. These patterns describe natural-language heuristics indicating syntax structures that suggest the presence of feature requests, for example: ‘‘[something] needs to be [verb]’’. They also used a parser that labels grammatical relationships between words within a sentence. This approach achieved fair precision and recall. The research group continued to develop the *User Request Referencer* (URR; Ciurumelea et al., 2017), which classifies reviews according to a custom taxonomy and then performs pattern-matching using Information Retrieval techniques to identify the affected source code artifacts. In Groen et al. (2017a), we theorized that online user feedback often contains redundancies. We obtained words specifically pertaining to *usability* from 360 user reviews and translated these into 16 linguistic patterns, achieving good precision—0.92 on average—but presumably limited

recall. Paech and Schneider (2020) and Schrieber et al. (2021) present research on the *User View Language*, an intermediate language with textual and visual notation of a software's outside view that aims to improve user-developer communication. Their work asserts that users have a different mental model than a system's developers, which will cause similar context to be described differently in utterances such as online user feedback. Their research suggests that the user view is quite consistent, regardless of technical background (Anders et al., 2022) or frequency of using an app (Anders et al., 2023). In spite of these works, we found that a comprehensive investigation of how suitable linguistic patterns are for user feedback analysis is missing. In RQ1, we used regular expressions to identify the end-user view on software quality.

Requirements classification. The proposed solutions for classifying problems into RE dimensions for structured or unstructured documents have most often employed NLP techniques backed by traditional ML models (for reviews, see, e.g., Santos et al., 2019a,b; Wang et al., 2019; Khan et al., 2019). This section focuses on research efforts of classifying dimensions of non-functional requirements or quality requirements, because this is the main focus of our work.

Before online user feedback was considered a viable source of information for RE, Cleland-Huang et al. (2006) used an Information Retrieval approach with a predefined fixed set of keywords extracted from standardized catalogs to detect and automatically classify NFRs in requirements specifications. This was extended with early demonstrations of using speech recognition to classify NFRs in meeting and interview recordings (Steele et al., 2006) and to analyze stakeholders' quality concerns from unstructured documents such as meeting minutes, interview notes, and memos (Cleland-Huang et al., 2007b), with the goal of reducing the number of overlooked NFRs through manual discovery. The work of Cleland-Huang et al. (2007b) bears the greatest resemblance to our work on RQ1. They created keyword taxonomies for *security* and *performance*, which they then used to iteratively retrain a classification algorithm to detect and classify (i.e., predict) NFRs through a series of experiments performed over 30 documents, increasing its recall from 62.6% to 79.9% and its precision from 14.7% to 20.7%. In RQ1, we detected and classified the quality concerns that crowd members express in online user feedback in a similar way.

NFRs are usually found in technical documents such as requirements specifications or SE contracts, which are typically well-curated and validated product requirement datasets. The analysis of such documents has become a popular field of study, and in recent years, techniques based predominantly on DL have been proposed to classify NFRs and other requirements (e.g., Navarro-Almanza et al., 2017; Tamai and Anzai, 2018; Hey et al., 2020; Sainani et al., 2020; Ajagbe and Zhao, 2022). The performance of these classifiers is most often validated on the *PROMISE* dataset (Sayyad Shirabad and Menzies, 2005; Cleland-Huang et al., 2007a) of user requirements. However, they face limitations when applied to online user feedback (Reddy Mekala et al., 2021; Dąbrowski et al., 2022). CrowdRE research has demonstrated the suitability of online user feedback for obtaining quality requirements (e.g., Groen et al., 2017a; Jha and Mahmoud, 2017; Lu and Liang, 2017), but unlike structured technical documents, online user feedback consists of unstructured, noisy, and inherently ambiguous texts not specifically intended for RE. It does not readily contain existing NFRs, which makes it more difficult to analyze. In their review, Dąbrowski et al. (2022) found that the analysis of NFRs comprises one of the three central use cases in online user feedback analysis. Through automated identification, classification, summarization, and sentiment analysis, NFRs can be elicited and problems can be identified, categorized into quality attributes, specified ad hoc with a documented rationale, and prioritized. Research works addressing the problem of (systematically) discovering quality requirements using classification methods have typically proposed methods employing traditional ML or DL (Zhou et al., 2014; Yang and Liang, 2015; Groen et al., 2017a; Lu and Liang, 2017; Jha and Mahmoud, 2017; Kurtanović and Maalej,

2017; Williams and Mahmoud, 2017; Stanik et al., 2019; Reddy Mekala et al., 2021). Although ML-based text mining tools in RE can achieve high recall values, their main weakness is usually that precision is often lacking (Jha and Mahmoud, 2017). This is why we investigated the suitability of other approaches for this classification problem.

6.3. LLMs in RE

In RQ3, we consider the suitability of LLMs using the example of ChatGPT to classify online user feedback according to RE-relevant dimensions. Due to the novelty of LLMs, little research on their application in RE had been published when we performed our investigation. A growing body of work in RE is considering the role of LLMs in education (Abdelfattah et al., 2023; Carvallo and Erazo-Garzón, 2023) and explores the generation of requirements (Brand, 2023; Bencheikh and Höglund, 2023), models (Ruan et al., 2023; Arulmohan et al., 2023; Bragilovski et al., 2024) or artifacts such as personas or interview scripts (Görer and Aydemir, 2023; Brand, 2023).

Regarding classification with the help of LLMs, like in RQ3, in the seminal work of Gilardi et al. (2023), ChatGPT using a zero-shot prompting strategy was found to have outperformed human raters on classifying tweets and news articles in a political science study. Classification was, among other things, by relevance, topic, sentiment, and stance. The only negative effect of the zero-shot prompting strategy might have been its diminished ability to determine the relevance of some texts due to a lack of examples. Gilardi et al. (2023) suggest that future work should consider few-shot prompting strategies, citing evidence that this strategy is best suited for LLMs (Brown et al., 2020). Although the domain is different, their work shows some similarities with our crowdsourced task in RQ2 and the use of this data in a study with LLMs in RQ3, which also classified the data according to relevance (Phases P1 & P2) and aspects that are similar to topics (Phases P3 & P4). In particular, their topic annotation task resembles our Phase 3' in terms of complexity.

So far, only few known works have investigated automated classification according to requirements-relevant content using LLMs as in our RQ3. El-Hajjaji et al. (2024) found that their ChatGPT models (3.5 Legacy [DaVinci], 3.5 Latest [Turbo], and 4) almost always outperformed LSTM- and SVM-based classifiers in binary classifications replicating those of Dalpiaz et al. (2019). Their conflicting results on the use of examples might suggest that these worsened the task location of the LLM (cf. Reynolds and McDonnell, 2021) or introduced bias. Recent works have also investigated the classification of requirements with learning. Wang et al. (2024) compared the ability of various language model types to predict requirements in two pre-existing annotated datasets and found DeBERTa to outperform Llama2 and RoBERTa. Ensembles were also found to achieve better results. van Can and Dalpiaz (2025) employed LLMs to identify and then classify requirements in backlog items in issue-tracking systems. They found that the encoder-only models they used, BERT and RoBERTa, outperformed the decoder-only models, ChatGPT 4 with and without in-context learning, Mistral 7B, and Llama 3.

At the time we started our research, little work specifically on prompt engineering existed or was known to us. We found helpful guidelines in Pattyn (2024) and El-Hajjaji et al. (2024), who suggested that the input data, context, procedural guidance classification schema, and expected structure of the output must be clearly described in the prompt syntax. In her practice report, Brand (2023) furthermore determined that instructions should be clear, specific, and written in a way that triggers the model to provide versatile responses. More useful resources are now available, including a catalog of prompt patterns (White et al., 2023), the *OpenPrompt* framework (Ding et al., 2022) for prompt learning and prompt management for research, the Automatic Prompt Engineer method (APE; Zhou et al., 2023) for automatically composing and selecting instructions, and datasets provided by Das et al. (2023), including one to evaluate zero-shot models through natural language inference.

7. Conclusion & future work

Identifying software product quality characteristics in inherently ambiguous online user feedback is a complex classification problem, consisting of multiclass and possibly multi-label classification according to some taxonomies such as ISO 25 010 (ISO, 2011). The automation potential for this classification problem is inhibited due to the lack of large training corpora, resulting in a low-data environment, which limits the potential for automated classifiers such as traditional ML or DL algorithms to be trained and perform this task with high accuracy. This is why in this research, we investigated and compared three low-data approaches in order to answer our main research question (MRQ): “Which low-data approach is most accurate in identifying quality aspects in online user feedback?”

We decomposed our MRQ into four research questions (RQs) to analyze and compare these approaches in their analysis of a sample of 1000 app store reviews. For RQ1, we analyzed traditional linguistic approaches through an adaptation of the NFR Method, through which we obtained a set of language patterns (LPs). Although it successfully led to a lexicon of keywords and the creation of a set of LPs, the vetting process was intensive work, and an experiment showed that the classification only predicted the quality characteristic of *reliability* sufficiently well. For RQ2, we extended prior research on the Kyōryoku method (van Vliet et al., 2020) to crowdsourcing the classification of quality aspects in online user feedback. An experiment showed that a series of smaller classification tasks achieved better results than a single, larger classification task. It demonstrated that it is possible for online user feedback to be classified into all quality aspects by a crowd of lay workers. For RQ3, we modeled an LLM pipeline after the structure used in RQ2 to determine its ability to achieve good results on this classification task with little training data. An experiment showed that ChatGPT 4o outperformed ChatGPT 4 Legacy in classifying by requirements relevance, and that using ChatGPT 4o with the instructions we provided to the crowd in Kyōryoku rather than the specifically engineered prompt somewhat boosted performance. This suggests that LLMs are capable of being used for this classification problem.

We performed a comparison of these three approaches to answer RQ4, which sought to identify which low-data approach is most suited for classifying quality aspects in online user feedback, considering effectiveness and trade-offs. Our results showed that although the use of crowdsourcing using the Kyōryoku method had a slight advantage over LLMs for classification by requirements relevance, the best configurations of both approaches proved to be equally suitable for classifying quality aspects, while the LP-based approach fell short. This suggests that both crowdsourcing and LLMs should be considered for *complex tasks* in RE, especially for purposes for which there is not enough data available to train ML classifiers. Specifically, the Kyōryoku method—introduced in van Vliet et al. (2020) and extended in this work—and our LLM pipeline with the associated artifacts can support the identification and classification of RE-relevant content in online user feedback. The implication of this finding is that the body of work on CrowdRE has been extended with two suitable approaches for this purpose.

In this work, we considered the three low-data approaches in isolation. A primary focus in future work should be to investigate whether a combination of approaches can help yield better results by compensating for each other's weaknesses and limitations. Ensemble methods for automatically classifying app reviews have previously been shown to perform better than individually applied techniques (Guzmán et al., 2015). We see several possibilities for how these approaches could support each other. The effective use of LLMs depends on their collaboration with human operators, for example, through supervision in learning. As a result, the best performance might be attained through a pipeline that, for example, uses automated analysis to verify the results from crowdsourcing and reattributes sentences for sequential phases. A keyword meta-model could help to expand the domain knowledge

of the LLM, and the output of a crowd or an LLM could be grouped through linguistic analyses to find inconsistencies in the heuristics applied when tagging, with the LLM in turn suggesting potential to relabel the data accordingly. Moreover, given that the three low-data approaches do not provide an algorithm that can be continuously trained to improve, we also see potential in the use of any of the approaches to construct a larger data corpus of online user feedback classified into quality aspects, which can, in turn, be used to train traditional ML and DL algorithms.

CRedit authorship contribution statement

Eduard C. Groen: Writing – original draft, Visualization, Validation, Methodology, Investigation, Formal analysis, Data curation, Conceptualization. **Fabiano Dalpiaz:** Writing – review & editing, Validation, Supervision, Methodology, Investigation, Conceptualization. **Martijn van Vliet:** Writing – original draft, Validation, Methodology, Investigation. **Boris Winter:** Writing – original draft, Validation, Investigation. **Joerg Doerr:** Writing – review & editing, Writing – original draft, Methodology, Conceptualization. **Sjaak Brinkkemper:** Writing – review & editing, Conceptualization.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgments

The authors would like to thank all the participants of the keywords & phrases elicitation workshops for RQ1, and Svenja Polst and Shivaji Yadav for their support in moderating the workshops. The authors thank Jonathan Ullrich and Dan Berry for their insightful advice on RQ3. This research did not receive specific grants from funding agencies in the public, commercial, or not-for-profit sectors.

Data availability

The online appendix with supplementary material, including additional content, experimental artifacts, and spreadsheets with data and results is available on Zenodo: <https://doi.org/10.5281/zenodo.15604749>.

References

- Abdelfattah, A.M., Ali, N.A., Elaziz, M.A., Ammar, H.H., 2023. Roadmap for software engineering education using ChatGPT. In: Proceedings of the 2nd International Conference on Artificial Intelligence Science and Applications in Industry and Society. CAISALS, IEEE, <http://dx.doi.org/10.1109/CAISALS59399.2023.10270477>, Article 14.
- Adepetu, A., Ahmed, K.A., Al Abd, Y., Al Zaabi, A., Svetinovic, D., 2012. CrowdREquire: A requirements engineering crowdsourcing platform. In: AAAI Spring Symposium: Wisdom of the Crowd. AAAI Technical Report SS-12-06, Article 1.
- Ajagbe, M., Zhao, L., 2022. Retraining a BERT model for transfer learning in requirements engineering: A preliminary study. In: Proceedings of the 30th International Requirements Engineering Conference. RE, IEEE, pp. 309–315. <http://dx.doi.org/10.1109/RE54965.2022.00046>.
- Alamer, G., Alyahya, S., 2023. A proposed approach to crowd selection in crowdsourced requirements engineering for mobile apps. In: Proceedings of the 7th International Conference on Information Systems Engineering. ICISE, ACM, pp. 1–5. <http://dx.doi.org/10.1145/3573926.3573927>.
- Ali, R., Solís, C., Omoronyia, I., Salehie, M., Nuseibeh, B., 2012. Social adaptation: When software gives users a voice. In: Proceedings of the 7th International Conference on Evaluation of Novel Approaches To Software Engineering. ENASE, Springer, pp. 75–84. <http://dx.doi.org/10.5220/0003991900750084>.
- Anders, M., Obaidi, M., Paech, B., Schneider, K., 2022. A study on the mental models of users concerning existing software. In: Gervasi, V., Vogelsang, A. (Eds.), Proceedings of Requirements Engineering – Foundation for Software Quality. REFSQ, In: LNCS 13216, Springer, pp. 256–271. http://dx.doi.org/10.1007/978-3-030-98464-9_18.

- Anders, M., Obaidi, M., Specht, A., Paech, B., 2023. What can be concluded from user feedback? – An empirical study. In: Proceedings of the 31st International Requirements Engineering Conference Workshops. REW, IEEE, pp. 122–128. <http://dx.doi.org/10.1109/REW57809.2023.00027>.
- Arulmohan, S., Meurs, M.-J., Mosser, S., 2023. Extracting domain models from textual requirements in the era of large language models. In: Proceedings of the 26th International Conference on Model Driven Engineering Languages and Systems Companion. MODELS-C, ACM/IEEE, pp. 580–587. <http://dx.doi.org/10.1109/MODELS-C59198.2023.00096>.
- Astegher, M., Busetta, P., Gabbasov, A., Pedrotti, M., Perini, A., Susi, A., 2023. Specifying requirements for collection and analysis of online user feedback. *Requir. Eng.* 28, 75–96. <http://dx.doi.org/10.1007/s00766-022-00387-3>.
- Bano, M., Zowghi, D., da Rimini, F., 2017. User satisfaction and system success: An empirical exploration of user involvement in software development. *Empir. Softw. Eng.* 22, 2339–2372. <http://dx.doi.org/10.1007/s10664-016-9465-1>.
- Bayerl, P.S., Paul, K.I., 2011. What determines inter-coder agreement in manual annotations? A meta-analytic investigation. *Comput. Linguist.* 37 (4), 699–725. http://dx.doi.org/10.1162/COLI_a.00074.
- Bencheikh, L., Höglund, N., 2023. Exploring the efficacy of ChatGPT in generating requirements: An experimental study. University of Gothenburg, Bachelor's thesis. <https://hdl.handle.net/2077/77957>.
- Berry, D.M., 2021. Empirical evaluation of tools for hairy requirements engineering tasks. *Empir. Softw. Eng.* 26 (6), 111. <http://dx.doi.org/10.1007/s10664-021-09986-0>.
- Bragilovski, M., van Can, A.T., Dalpiaz, F., Sturm, A., 2024. Deriving domain models from user stories: Human vs. Machines. In: Proceedings of the 32nd International Requirements Engineering Conference. RE, IEEE, pp. 31–42. <http://dx.doi.org/10.1109/RE59067.2024.00014>.
- Brand, A., 2023. Wie viel RE geht mit generativer KI? Ein Selbstversuch. Presentation at Summit Community Days Requirements Engineering, Leipzig, Germany, 21 November 2023.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., et al., 2020. Language models are few-shot learners. In: Larochelle, H., Ranzato, M., Hadsell, R., et al. (Eds.), *Advances in Neural Information Processing Systems*. Vol. 33, Curran Associates, pp. 1877–1901. https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfb4967418bfb8ac142f64a-Paper.pdf.
- Campbell, D.T., Stanley, J.C., 1963. *Experimental and quasi-experimental designs for research*. Houghton Mifflin, Boston, MA, USA.
- van Can, A.T., Dalpiaz, F., 2025. Locating requirements in backlog items: Content analysis and experiments with large language models. *Inf. Softw. Technol.* 179, 107644. <http://dx.doi.org/10.1016/j.infsof.2024.107644>.
- Carvalho, J.P., Erazo-Garzon, L., 2023. On the use of ChatGPT to support requirements engineering teaching and learning process. In: Berrezueta, S. (Ed.), *Proceedings of the 18th Latin American Conference on Learning Technologies*. LACLO, Springer Nature, Singapore, pp. 328–342. http://dx.doi.org/10.1007/978-981-99-7353-8_25.
- Castro-Herrera, C., Cleland-Huang, J., Mobasher, B., 2009. Enhancing stakeholder profiles to improve recommendations in online requirements elicitation. In: *Proceedings of the 17th International Requirements Engineering Conference*. RE, IEEE, pp. 37–46. <http://dx.doi.org/10.1109/RE.2009.20>.
- Ceararu, I., Lazar, J., Bessiere, K., Robinson, J., Shneiderman, B., 2004. Determining causes and severity of end-user frustration. *Int. J. Human-Comput. Interact.* 17 (3), 333–356. <http://dx.doi.org/10.1207/s15327590ijhc1703.3>.
- Chung, L., Nixon, B.A., Yu, E., Mylopoulos, J., 2012. Non-functional requirements in software engineering. In: *International Series in Software Engineering*, vol. 5, Springer Science & Business Media, <http://dx.doi.org/10.1007/978-1-4615-5269-7>.
- Chung, L., do Prado Leite, J.C.S., 2009. On non-functional requirements in software engineering. In: *Conceptual Modeling: Foundations and Applications*. In: LNCS 5600, Springer, pp. 363–379. http://dx.doi.org/10.1007/978-3-642-02463-4_19.
- Ciurumelea, A., Schaufelbühl, A., Panichella, S., Gall, H.C., 2017. Analyzing reviews and code of mobile apps for better release planning. In: *Proceedings of the 24th International Conference on Software Analysis, Evolution and Reengineering*. SANER, IEEE, pp. 91–102. <http://dx.doi.org/10.1109/SANER.2017.7884612>.
- Cleland-Huang, J., Mazrouee, S., Huang, L., Port, D., 2007a. PROMISE-nfr. Zenodo, record 268542, <http://dx.doi.org/10.5281/zenodo.268542>, [Data set].
- Cleland-Huang, J., Settimi, R., Zou, X., Solc, P., 2006. The detection and classification of non-functional requirements with application to early aspects. In: *Proceedings of the 14th International Requirements Engineering Conference*. RE, IEEE, pp. 39–48. <http://dx.doi.org/10.1109/RE.2006.65>.
- Cleland-Huang, J., Settimi, R., Zou, X., Solc, P., 2007b. Automated classification of non-functional requirements. *Requir. Eng.* 12 (2), 103–120. <http://dx.doi.org/10.1007/s00766-007-0045-1>.
- Cysneiros, L.M., do Prado Leite, J.C.S., de Melo Sabat Neto, J., 2001. A framework for integrating non-functional requirements into conceptual models. *Requir. Eng.* 6 (2), 97–115. <http://dx.doi.org/10.1007/s007660170008>.
- Dąbrowski, J., Letier, E., Perini, A., Susi, A., 2022. Mining user feedback for software engineering: Use cases and reference architecture. In: *Proceedings of the 30th International Requirements Engineering Conference Workshops*. REW, IEEE, pp. 114–126. <http://dx.doi.org/10.1109/REW55302.2022.00023>.
- Dalpiaz, F., Dell'Anna, D., Aydemir, F.B., Çevikol, S., 2019. Requirements classification with interpretable machine learning and dependency parsing. In: *Proceedings of the 27th International Requirements Engineering Conference*. RE, IEEE, pp. 142–152. <http://dx.doi.org/10.1109/RE.2019.00025>.
- Das, S., Deb, N., Cortesi, A., Chaki, N., 2023. Zero-shot learning for named entity recognition in software specification documents. In: *Proceedings of the 31st International Requirements Engineering Conference*. RE, IEEE, pp. 100–110. <http://dx.doi.org/10.1109/RE57278.2023.00019>.
- Delima, R., Riastiawan, M., Ashari, A., 2023. Design of automatic user identification framework in crowdsourcing requirements engineering: User mapping and system architecture. *Creative Commun. Innov. Technol.* 16 (1), 54–67. <http://dx.doi.org/10.33050/ccit.v16i1>.
- Dell'Anna, D., Aydemir, F.B., Dalpiaz, F., 2023. Evaluating classifiers in SE research: The ECSE pipeline and two replication studies. *Empir. Softw. Eng.* 28 (3), <http://dx.doi.org/10.1007/s10664-022-10243-1>, Article 3.
- Demšar, J., 2006. Statistical comparisons of classifiers over multiple data sets. *J. Mach. Learn. Res.* 7, 1–30. <http://dx.doi.org/10.5555/1248547.1248548>.
- Dermeval, D., Vilela, J., Bittencourt, I.I., Castro, J., Isotani, S., Brito, P.H.d.S., Silva, A., 2016. Applications of ontologies in requirements engineering: A systematic review of the literature. *Requir. Eng.* 21 (4), 405–437. <http://dx.doi.org/10.1007/s00766-015-0222-6>.
- Deshpande, G., Motger, Q., Palomares, C., Kamra, I., Biesalska, K., Franch, X., Ruhe, G., Ho, J., 2020. Requirements dependency extraction by integrating active learning with ontology-based retrieval. In: *Proceedings of the 28th International Requirements Engineering Conference*. RE, IEEE, pp. 78–89. <http://dx.doi.org/10.1109/RE48521.2020.00020>.
- Ding, N., Hu, S., Zhao, W., Chen, Y., Liu, Z., Zheng, H., Sun, M., 2022. Open-Prompt: An open-source framework for prompt-learning. In: *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*. Association for Computational Linguistics, pp. 105–113. <http://dx.doi.org/10.18653/v1/2022.acl-demo.10>.
- Dörr, J., 2011. Elicitation of a complete set of non-functional requirements (Ph.D. thesis). University of Kaiserslautern, <http://dx.doi.org/10.24406/publica-fhg-278948>, PhD Theses in Experimental Software Engineering, Vol. 34. Stuttgart, Germany: Fraunhofer Verlag.
- El-Hajjani, A., Fafin, N., Salinesi, C., 2024. Which AI technique is better to classify requirements? An experiment with SVM, LSTM, and ChatGPT. In: *Requirements Engineering – Foundation for Software Quality (REFSQ) Joint Proceedings of the Co-located Events*. CEUR Workshop Proceedings 3672, NLP4RE article 4. <https://ceur-ws.org/Vol-3672/NLP4RE-paper2.pdf>.
- Fawcett, T., 2006. An introduction to ROC analysis. *Pattern Recognit. Lett.* 27 (8), 861–874. <http://dx.doi.org/10.1016/j.patrec.2005.10.010>.
- Ferrari, A., Ginde, G., 2025. *Handbook on natural language processing for requirements engineering*. Springer, Cham, Switzerland, http://dx.doi.org/10.1007/978-3-031-73143-3_1.
- Fitzgerald, B., 2018. Crowdsourcing software development: Silver bullet or lead balloon. In: *Proceedings of the 5th International Workshop on Artificial Intelligence for Requirements Engineering*. AIRE, IEEE, pp. 29–30. <http://dx.doi.org/10.1109/AIRE.2018.00010>.
- Fu, B., Lin, J., Lei, L., Faloutsos, C., Hong, J.I., Sadeh, N.M., 2013. Why people hate your app: Making sense of user feedback in a mobile app store. In: *Proceedings of the 19th International Conference on Knowledge Discovery and Data Mining*. ACM SIGKDD, pp. 1276–1284. <http://dx.doi.org/10.1145/2487575.2488202>.
- Ghezzi, A., Gabelloni, D., Martini, A., Natalicchio, A., 2018. Crowdsourcing: A review and suggestions for future research. *Int. J. Manag. Rev.* 20 (2), 343–363. <http://dx.doi.org/10.1111/ijmr.12135>.
- Gilardi, F., Alizadeh, M., Kubli, M., 2023. ChatGPT outperforms crowd workers for text-annotation tasks. *Proc. Natl. Acad. Sci.* 120 (30), e2305016120. <http://dx.doi.org/10.1073/pnas.2305016120>.
- Glinz, M., 2007. On non-functional requirements. In: *Proceedings of the 15th International Requirements Engineering Conference*. RE, IEEE, pp. 21–26. <http://dx.doi.org/10.1109/RE.2007.45>.
- Glinz, M., 2019. CrowdRE: Achievements, opportunities and pitfalls. In: *Proceedings of the 27th International Requirements Engineering Conference Workshops*. REW, IEEE, pp. 172–173. <http://dx.doi.org/10.1109/REW.2019.00036>.
- Glinz, M., Seyff, N., Bühne, S., Franch, X., Lauenroth, K., 2023. Towards a modern quality framework. In: Schneider, K., Dalpiaz, F., Horkoff, J. (Eds.), *Proceedings of the 31st International Requirements Engineering Conference Workshops*. REW, IEEE, pp. 357–361. <http://dx.doi.org/10.1109/REW57809.2023.00067>.
- Görer, B., Aydemir, F.B., 2023. Generating requirements elicitation interview scripts with large language models. In: *Proceedings of the 31st International Requirements Engineering Conference Workshops*. REW, IEEE, pp. 44–51. <http://dx.doi.org/10.1109/REW57809.2023.00015>.
- Groen, E.C., 2015. Crowd out the competition: Gaining market advantage through crowd-based requirements engineering. In: *Proceedings of the 1st International Workshop on Crowd-Based Requirements Engineering*. CrowdRE, IEEE, pp. 13–18. <http://dx.doi.org/10.1109/CrowdRE.2015.7367583>.
- Groen, E.C., Dalpiaz, F., van Vliet, M., Winter, B., Doerr, J., Brinkkemper, S., 2025. Classification of quality characteristics in online user feedback using linguistic analysis, crowdsourcing and LLMs: Online appendix and supplementary material. Zenodo, record 15604749, <https://doi.org/10.5281/zenodo.15604749>.

- Groen, E.C., Doerr, J., Adam, S., 2015. Towards crowd-based requirements engineering: A research preview. In: Fricker, S.A., Schneider, K. (Eds.), *Proceedings of Requirements Engineering – Foundation for Software Quality. REFSQ*, In: LNCS 9013, Springer, pp. 247–253. http://dx.doi.org/10.1007/978-3-319-16101-3_16.
- Groen, E.C., Kocpzyńska, S., Hauer, M.P., Krafft, T.D., Doerr, J., 2017a. Users — the hidden software product quality experts? A study on how app users report quality aspects in online reviews. In: *Proceedings of the 25th International Requirements Engineering Conference. RE, IEEE*, pp. 80–89. <http://dx.doi.org/10.1109/RE.2017.73>.
- Groen, E.C., Ochs, M., 2019. CrowdRE, user feedback and GDPR: Towards tackling GDPR implications with adequate technical and organizational measures in an effort-minimal way. In: *Proceedings of the 27th International Requirements Engineering Conference Workshops. REW, IEEE*, pp. 180–185. <http://dx.doi.org/10.1109/REW.2019.00038>.
- Groen, E.C., Schowalter, J., Kocpzyńska, S., Polst, S., Alvani, S., 2018. Is there really a need for using NLP to elicit requirements? A benchmarking study to assess scalability of manual analysis. In: Schmid, K., Spoletini, P. (Eds.), *Requirements Engineering – Foundation for Software Quality (REFSQ) Joint Proceedings of the Co-Located Events. CEUR Workshop Proceedings 2075, NLP4RE article 11*. https://ceur-ws.org/Vol-2075/NLP4RE_paper11.pdf.
- Groen, E.C., Seyff, N., Ali, R., Dalpiaz, F., Doerr, J., Guzmán, E., Hosseini, M., Marco, J., Oriol, M., Perini, A., Stade, M., 2017b. The crowd in requirements engineering: The landscape and challenges. *IEEE Softw.* 34 (2), 44–52. <http://dx.doi.org/10.1109/MS.2017.33>.
- Guzmán, E., El-Haliby, M., Brügge, B., 2015. Ensemble methods for app review classification: An approach for software evolution (N). In: *Proceedings of the 30th International Conference on Automated Software Engineering. ASE, IEEE/ACM*, pp. 771–776. <http://dx.doi.org/10.1109/ASE.2015.88>.
- Hertzum, M., Hornbæk, K., 2023. Frustration: Still a common user experience. *ACM Trans. Comput.-Hum. Interact.* 30 (3), <http://dx.doi.org/10.1145/3582432>, Article 42.
- Hey, T., Keim, J., Koziol, A., Tichy, W.F., 2020. NoBERT: Transfer learning for requirements classification. In: *Proceedings of the 28th International Requirements Engineering Conference. RE, IEEE*, pp. 169–179. <http://dx.doi.org/10.1109/RE48521.2020.00028>.
- Hosmer, D.W.J., Lemeshow, S., Sturdivant, R.X., 2013. *Applied logistic regression*, third ed. Wiley, Hoboken, NJ. <http://dx.doi.org/10.1002/9781118548387>.
- Hosseini, M., Groen, E.C., Shahri, A., Ali, R., 2017. CRAFT: A crowd-annotated feedback technique. In: *Proceedings of the 2nd International Workshop on Crowd-Based Requirements Engineering. CrowdRE, IEEE*, pp. 170–175. <http://dx.doi.org/10.1109/REW.2017.27>.
- Hosseini, M., Phalp, K.T., Taylor, J., Ali, R., 2014a. The four pillars of crowdsourcing: A reference model. In: *Proceedings of the 8th International Conference on Research Challenges in Information Science. RCIS, IEEE*, pp. 1–12. <http://dx.doi.org/10.1109/RCIS.2014.6861072>.
- Hosseini, M., Phalp, K.T., Taylor, J., Ali, R., 2014b. Towards crowdsourcing for requirements engineering. In: *Requirements Engineering – Foundation for Software Quality (REFSQ) Joint Proceedings of the Co-Located Events. CEUR Workshop Proceedings 1138, Empirical Track article 2*. <https://ceur-ws.org/Vol-1138/et2.pdf>.
- Hosseini, M., Shahri, A., Phalp, K.T., Taylor, J., Ali, R., 2015a. Crowdsourcing: A taxonomy and systematic mapping study. *Comput. Sci. Rev.* 17 (C), 43–69. <http://dx.doi.org/10.1016/j.cosrev.2015.05.001>.
- Hosseini, M., Shahri, A., Phalp, K.T., Taylor, J., Ali, R., Dalpiaz, F., 2015b. Configuring crowdsourcing for requirements elicitation. In: *Proceedings of the 9th International Conference on Research Challenges in Information Science. RCIS, IEEE*, pp. 133–138. <http://dx.doi.org/10.1109/RCIS.2015.7128873>.
- Howe, J., 2006. The rise of crowdsourcing. *Wired* 14 (6), 1–4. <https://www.wired.com/2006/06/crowds/>.
- Howe, J., 2010. Crowdsourcing: A definition. <https://crowdsourcing.typepad.com/>. (Accessed 15 June 2025).
- Iacob, C., Harrison, R., 2013. Retrieving and analyzing mobile apps feature requests from online reviews. In: *Proceedings of the 10th Working Conference on Mining Software Repositories. MSR, IEEE*, pp. 41–44. <http://dx.doi.org/10.1109/MSR.2013.6624001>.
- ISO, 2011. *Systems and software engineering – Systems and software quality requirements and evaluation (SQuaRE) – System and software quality models*. Standard ISO/IEC 25010:2011, International Organization for Standardization, Geneva, Switzerland. <https://www.iso.org/standard/35733.html>.
- ISO/IEC, 2023. *ISO/IEC 25010 – Systems and software engineering – Systems and software quality requirements and evaluation (SQuaRE) – Product quality model*. Standard ISO/IEC 25010:2011, International Organization for Standardization, Geneva, Switzerland. <https://www.iso.org/standard/78176.html>.
- Jha, N., Mahmoud, A., 2017. Mining user requirements from application store reviews using frame semantics. In: Grünbacher, P., Perini, A. (Eds.), *Proceedings of Requirements Engineering – Foundation for Software Quality. REFSQ*, In: LNCS 10153, Springer, pp. 237–287. http://dx.doi.org/10.1007/978-3-319-54045-0_20.
- Jhangiani, R.S., Chiang, I.-C.A., Cuttler, C., Leighton, D.C., 2019. Research methods in psychology, fourth ed. Kwantlen Polytechnic University, Fairmount, Montreal, Canada. <http://dx.doi.org/10.17605/OSF.IO/HF7DQ>.
- Jurafsky, D., Martin, J.H., 2009. *Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition*, second ed. In: *Prentice Hall Series in Artificial Intelligence*, Prentice Hall, Pearson Education International, Upper Saddle River, NJ, USA.
- Khan, J.A., Liu, L., Wen, L., Ali, R., 2019. Crowd intelligence in requirements engineering: Current status and future directions. In: Knauss, E., Goedicke, M. (Eds.), *Proceedings of Requirements Engineering – Foundation for Software Quality. REFSQ*, In: LNCS 11412, Springer, pp. 245–261. http://dx.doi.org/10.1007/978-3-030-15538-4_18.
- Khan, H.H., Malik, M.N., Alotaibi, Y., Alsufyani, A., Alghamdi, S., 2021. Crowdsourced requirements engineering challenges and solutions: A software industry perspective. *Comput. Syst. Sci. Eng.* 39 (2), 221–236. <http://dx.doi.org/10.32604/csse.2021.016510>.
- Kittur, A., Smus, B., Khamkar, S., Kraut, R.E., 2011. CrowdForge: Crowdsourcing complex work. In: *Proceedings of the 24th Annual ACM Symposium on User Interface Software and Technology. UIST, ACM*, pp. 43–52. <http://dx.doi.org/10.1145/2047196.2047202>.
- Kurtanović, Z., Maalej, W., 2017. Automatically classifying functional and non-functional requirements using supervised machine learning. In: *Proceedings of the 25th International Requirements Engineering Conference. RE, IEEE*, pp. 490–495. <http://dx.doi.org/10.1109/RE.2017.82>.
- Kutner, M.H., Nachtsheim, C.J., Neter, J., Li, W., 2005. *Applied linear statistical models*, fifth ed. McGraw-Hill/Irwin, Chicago, IL, and Boston, MA.
- van Lamsweerde, A., 2001. Goal-oriented requirements engineering: A guided tour. In: *Proceedings of the 5th International Symposium on Requirements Engineering. RE, IEEE*, pp. 249–262. <http://dx.doi.org/10.1109/ISRE.2001.948567>.
- Lim, S.L., Quercia, D., Finkelstein, A., 2010. StakeSource: Harnessing the power of crowdsourcing and social networks in stakeholder analysis. In: *Proceedings of the 32nd International Conference on Software Engineering. ICSE, Vol. 2, ACM/IEEE*, pp. 239–242. <http://dx.doi.org/10.1145/1810295.1810340>.
- Lu, M., Liang, P., 2017. Automatic classification of non-functional requirements from augmented app user reviews. In: *Proceedings of the 21st International Conference on Evaluation and Assessment in Software Engineering. EASE*, pp. 344–353. <http://dx.doi.org/10.1145/3084226.3084241>.
- Mao, K., Capra, L., Harman, M., Jia, Y., 2017. A survey of the use of crowdsourcing in software engineering. *J. Syst. Softw.* 126, 57–84. <http://dx.doi.org/10.1016/j.jss.2016.09.015>.
- Mardia, K.V., 1970. Measures of multivariate skewness and kurtosis with applications. *Biom.* 57 (3), 519–530. <http://dx.doi.org/10.1093/biomet/57.3.519>.
- Mizrahi, M., Kaplan, G., Malkin, D., Dror, R., Shahaf, D., Stanovsky, G., 2024. State of what art? A call for multi-prompt LLM evaluation. *Trans. Assoc. Comput. Linguist.* 12, 933–949. http://dx.doi.org/10.1162/tacl_a.00681.
- Motger, Q., Borrell, R., Palomares, C., Marco, J., 2019. OpenReq-DD: A requirements dependency detection tool. In: *Requirements Engineering – Foundation for Software Quality (REFSQ) Joint Proceedings of the Co-Located Events. CEUR Workshop Proceedings 2376, Tool Demonstrations article 2*. https://ceur-ws.org/Vol-2376/NLP4RE19_paper01.pdf.
- Navarro-Almanza, R., Juárez-Ramírez, R., Licea, G., 2017. Towards supporting software engineering using deep learning: A case of software requirements classification. In: *Proceedings of the 5th International Conference in Software Engineering Research and Innovation. CONISOF, IEEE*, pp. 116–120. <http://dx.doi.org/10.1109/CONISOF.2017.00021>.
- Oluwatofunmi, A., Ebenezer, O., Nzechukwu, O., 2020. Crowd requirement rating technique (CrowdReRaT) model for crowd sourcing. *Int. J. Comput. Appl.* 176 (22), 9–14. <http://dx.doi.org/10.5120/ijca2020920178>.
- Paech, B., Schneider, K., 2020. How do users talk about software? Searching for common ground. In: *Proceedings of the 1st Workshop on Ethics in Requirements Engineering Research and Practice. REthics, IEEE*, pp. 11–14. <http://dx.doi.org/10.1109/REthics51204.2020.00008>.
- Pagano, D., Maalej, W., 2013. User feedback in the appstore: An empirical study. In: *Proceedings of the 21st International Requirements Engineering Conference. RE, IEEE*, pp. 125–134. <http://dx.doi.org/10.1109/RE.2013.6636712>.
- Panichella, S., Di Sorbo, A., Guzmán, E., Visaggio, C.A., Canfora, G., Gall, H.C., 2015. How can I improve my app? Classifying user reviews for software maintenance and evolution. In: *Proceedings of the 31st IEEE International Conference on Software Maintenance and Evolution. ICSME, IEEE*, pp. 281–290. <http://dx.doi.org/10.1109/ICSM.2015.7332474>.
- Patryn, F., 2024. The value of generative AI for qualitative research: A pilot study. *J. Data Sci. Intell. Syst. Online First*, <http://dx.doi.org/10.47852/bonviewJDSIS42022964>.
- Pohl, K., Rupp, C., 2015. *Requirements engineering fundamentals: A study guide for the Certified Professional for Requirements Engineering Exam – Foundation Level – IREB Compliant*, second ed. Rocky Nook, Santa Barbara, CA, USA.
- Posch, L., Bleier, A., Flöck, F., Lechner, C.M., Kinder-Kurlanda, K., Helic, D., Strohmaier, M., 2022. Characterizing the global crowd workforce: A cross-country comparison of crowdworker demographics. *Hum. Comput.* 9 (1), 22–57. <http://dx.doi.org/10.15346/hc.v9i1.106>.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P.J., 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.* 21 (140), 1–67. <http://jmlr.org/papers/v21/20-074.html>.

- Reddy Mekala, R., Irfan, A., Groen, E.C., Porter, A., Lindvall, M., 2021. Classifying user requirements from online feedback in small dataset environments using deep learning. In: Proceedings of the 29th International Requirements Engineering Conference. RE, IEEE, pp. 139–149. <http://dx.doi.org/10.1109/RE51729.2021.00020>.
- Renzel, D., Behrendt, M., Klamma, R., Jarke, M., 2013. Requirements Bazaar: Social requirements engineering for community-driven innovation. In: Proceedings of the 21st International Requirements Engineering Conference. RE, IEEE, pp. 326–327. <http://dx.doi.org/10.1109/RE.2013.6636738>.
- Retelny, D., Robaszkiewicz, S., To, A., Lasecki, W.S., Patel, J., Rahmati, N., Doshi, T., Valentine, M., Bernstein, M.S., 2014. Expert crowdsourcing with flash teams. In: Proceedings of the 27th Annual Symposium on User Interface Software and Technology. UIST, ACM, pp. 75–85. <http://dx.doi.org/10.1145/2642918.2647409>.
- Reynolds, L., McDonnell, K., 2021. Prompt programming for large language models: Beyond the few-shot paradigm. In: Extended Abstracts of the Conference on Human Factors in Computing Systems CHI-EA. ACM, pp. 314:1–314:7. <http://dx.doi.org/10.1145/3411763.3451760>.
- Ruan, K., Chen, X., Jin, Z., 2023. Requirements modeling aided by ChatGPT: An experience in embedded systems. In: Proceedings of the 31st International Requirements Engineering Conference Workshops. REW, IEEE, pp. 170–177. <http://dx.doi.org/10.1109/REW57809.2023.00035>.
- Sachs, J., 2005. *The End of Poverty: Economic Possibilities for Our Time*. Penguin Press.
- Sainani, A., Anish, P.R., Joshi, V., Ghaisas, S., 2020. Extracting and classifying requirements from software engineering contracts. In: Proceedings of the 28th International Requirements Engineering Conference. RE, IEEE, pp. 147–157. <http://dx.doi.org/10.1109/RE48521.2020.00026>.
- Santos, R., Groen, E.C., Villela, K., 2019a. An overview of user feedback classification approaches. In: Requirements Engineering – Foundation for Software Quality (REFSQ) Joint Proceedings of the Co-Located Events. CEUR Workshop Proceedings 2376, NLP4RE article 7. https://ceur-ws.org/Vol-2376/NLP4RE19_paper11.pdf.
- Santos, R., Groen, E.C., Villela, K., 2019b. A taxonomy for user feedback classifications. In: Requirements Engineering – Foundation for Software Quality (REFSQ) Joint Proceedings of the Co-Located Events. CEUR Workshop Proceedings 2376, NLP4RE article 5. https://ceur-ws.org/Vol-2376/NLP4RE19_paper10.pdf.
- Sari, D.A.P., Putri, A.Y., Hanggareni, M., Anjani, A., Siswondo, M.L.O., Raharjana, I.K., 2021. Crowdsourcing as a tool to elicit software requirements. In: Alfinyah, C., Fatmawati, Windarto (Eds.), International Conference on Mathematics, Computational Sciences and Statistics (ICoMCoS) 2020. In: AIP Conference Proceedings, vol. 2329, AIP, Article 050001, <https://doi.org/10.1063/5.0042134>.
- Sayyad Shirabad, J., Menzies, T.J., 2005. The PROMISE repository of software engineering databases. School of Information Technology and Engineering, University of Ottawa, Canada, <http://promise.site.uottawa.ca/SERepository>.
- Schenk, E., Guittard, C., 2011. Towards a characterization of crowdsourcing practices. J. Innov. Econ. Manag. 1 (7), 93–107. <http://dx.doi.org/10.3917/jie.007.0093>.
- Schrieber, H., Anders, M., Paech, B., Schneider, K., 2021. A vision of understanding the users' view on software. In: Requirements Engineering – Foundation for Software Quality (REFSQ) Joint Proceedings of the Co-Located Events. CEUR Workshop Proceedings 2857, NLP4RE article 3. <https://ceur-ws.org/Vol-2857/nlp4re2.pdf>.
- Sihag, M., Li, Z.S., Dash, A., Arony, N.N., Devathanan, K., Ernst, N., Albu, A.B., Damian, D., 2023. A data-driven approach for finding requirements relevant feedback from TikTok and YouTube. In: Proceedings of the 31st International Requirements Engineering Conference. RE, IEEE, pp. 111–122. <http://dx.doi.org/10.1109/RE57278.2023.00020>.
- Snijders, R., Dalpiaz, F., Brinkkemper, S., Hosseini, M., Ali, R., Ozum, A., 2015. Refine: A gamified platform for participatory requirements engineering. In: Proceedings of the 1st International Workshop on Crowd-Based Requirements Engineering. CrowdRE, IEEE, pp. 1–6. <http://dx.doi.org/10.1109/CrowdRE.2015.7367581>.
- Stanik, C., Häring, M., Maalej, W., 2019. Classifying multilingual user feedback using traditional machine learning and deep learning. In: Proceedings of the 27th International Requirements Engineering Conference Workshops. REW, IEEE, pp. 220–226. <http://dx.doi.org/10.1109/REW.2019.00046>.
- Steele, A., Arnold, J., Cleland-Huang, J., 2006. Speech detection of stakeholders' non-functional requirements. In: Proceedings of the 1st International Workshop on Multimedia Requirements Engineering. MERE, ACM, Article 2. <https://doi.org/10.1109/MERE.2006.5>.
- Stol, K.-J., Fitzgerald, B., 2014. Two's company, three's a crowd: A case study of crowdsourcing software development. In: Proceedings of the 36th International Conference on Software Engineering. ICSE, pp. 187–198, Technical Track. <https://doi.org/10.1145/2568225.2568249>.
- Tamai, T., Anzai, T., 2018. Quality requirements analysis with machine learning. In: Proceedings of the 13th International Conference on Evaluation of Novel Approaches To Software Engineering. ENASE, SciTePress, pp. 241–248. <http://dx.doi.org/10.5220/0006694502410248>.
- Tizard, J., Rietz, T., Blincoe, K., 2020. Voice of the users: A demographic study of software feedback behaviour. In: Proceedings of the 28th International Requirements Engineering Conference. RE, IEEE, pp. 55–65. <http://dx.doi.org/10.1109/RE48521.2020.00018>.
- Tukey, J.W., 1949. Comparing individual means in the analysis of variance. Biom. 5 (2), 99–114. <http://dx.doi.org/10.2307/3001913>.
- Valentine, M.A., Retelny, D., To, A., Rahmati, N., Doshi, T., Bernstein, M.S., 2017. Flash organizations: Crowdsourcing complex work by structuring crowds as organizations. In: Proceedings of the Conference on Human Factors in Computing Systems. CHI, ACM, pp. 3523–3537. <http://dx.doi.org/10.1145/3025453.3025811>.
- van Vliet, M., Groen, E.C., Dalpiaz, F., Brinkkemper, S., 2020. Identifying and classifying user requirements in online feedback via crowdsourcing. In: Madhavji, N., Pasquale, L., Ferrari, A., Gnesi, S. (Eds.), Proceedings of Requirements Engineering – Foundation for Software Quality. REFSQ, In: LNCS 12045, Springer, pp. 143–159. http://dx.doi.org/10.1007/978-3-030-44429-7_11.
- Vogelsang, A., 2024. Prompting the future: Integrating generative LLMs and requirements engineering. In: Requirements Engineering – Foundation for Software Quality (REFSQ) Joint Proceedings of the Co-Located Events. CEUR Workshop Proceedings 3672, NLP4RE article 2. <https://ceur-ws.org/Vol-3672/NLP4RE-keynote1.pdf>.
- Wang, C., Daneva, M., van Sinderen, M., Liang, P., 2019. A systematic mapping study on crowdsourced requirements engineering using user feedback. J. Softw.: Evol. Process. 31 (10), e2199. <http://dx.doi.org/10.1002/smr.2199>.
- Wang, K., Zhang, F., Sabetzadeh, M., 2024. Automated requirements demarcation using large language models: An empirical study. In: Requirements Engineering – Foundation for Software Quality (REFSQ) Joint Proceedings of the Co-Located Events. CEUR Workshop Proceedings 3672, NLP4RE article 5. <https://ceur-ws.org/Vol-3672/NLP4RE-paper4.pdf>.
- White, J., Fu, Q., Hays, S., Sandborn, M., Olea, C., Gilbert, H., Elnashar, A., Spencer-Smith, J., Schmidt, D.C., 2023. A prompt pattern catalog to enhance prompt engineering with ChatGPT. arXiv:2302.11382 [cs.SE]. <https://doi.org/10.48550/arXiv.2302.11382>.
- Wieringa, R.J., 2014. Design science methodology for information systems and software engineering. Springer, Berlin, Heidelberg, <http://dx.doi.org/10.1007/978-3-662-43839-8>.
- Williams, G., Mahmoud, A., 2017. Mining Twitter feeds for software user requirements. In: Proceedings of the 25th International Requirements Engineering Conference. RE, IEEE, pp. 1–10. <http://dx.doi.org/10.1109/RE.2017.14>.
- Wohlin, C., Runeson, P., Höst, M., Ohlsson, M.C., Regnell, B., Wesslén, A., 2000. Empirical strategies. In: Experimentation in Software Engineering: An Introduction. Springer US, Boston, MA, USA, pp. 7–24. http://dx.doi.org/10.1007/978-1-4615-4625-2_2.
- Wouters, J., Menkveld, A., Brinkkemper, S., Dalpiaz, F., 2022. Crowd-based requirements elicitation via pull feedback: Method and case studies. Requir. Eng. 27, 429–455. <http://dx.doi.org/10.1007/s00766-022-00384-6>.
- Yang, H., Liang, P., 2015. Identification and classification of requirements from app user reviews. In: Proceedings of the 15th International Conference on Software Engineering and Knowledge Engineering. SEKE, Article 63. <https://doi.org/10.18293/SEKE2015-063>.
- Zhou, Y., Muresanu, A.I., Han, Z., Paster, K., Pitis, S., Chan, H., Ba, J., 2023. Large language models are human-level prompt engineers. arXiv:2211.01910 [cs.LG]. <https://doi.org/10.48550/arXiv.2211.01910>.
- Zhou, Y., Tong, Y., Gu, R., Gall, H., 2014. Combining text mining and data mining for bug report classification. In: Proceedings of the 30th International Conference on Software Maintenance and Evolution. ICSME, IEEE, pp. 311–320. <http://dx.doi.org/10.1109/ICSME.2014.53>.