



Utrecht
University

Enhancing the ECSEER pipeline: Evaluating Large Language Model Classifiers in SE Research

Ruben van der Zeijden (7081111)

Supervisor: Dr. F.B. Aydemir

Second Reader: Dr. D. Dell'Anna

Abstract

The various research papers in the field of Software Engineering (SE) that use classification algorithms, LLMs, or other machine learning methods to obtain their results differ in how many and which evaluation metrics are reported, whether or not significance tests are performed, and which steps are taken to aid reproducibility. The ECSEER pipeline was designed to mitigate this issue by providing a step-by-step pipeline that researchers can use to report their results when classification algorithms are used, and the pipeline was empirically shown to be effective in replicating the findings of several studies, as well as producing additional findings and occasionally contradicting the findings of the original papers. However, the ECSEER pipeline is designed for evaluating simple classifiers and gives no specific recommendations for LLM classifiers, despite LLMs being an increasingly popular choice for classification tasks. The goal of this thesis is to expand the ECSEER pipeline for the use of LLMs by adding recommendations from LLM4SE research and related fields. First, an exploratory mapping study will be done of SE studies released in 2024 that use LLM classifiers, summarising which steps are taken and which metrics are (or aren't) reported, in order to get an overview of the current state of LLM4SE research. Subsequently, we design an enhanced version of the ECSEER pipeline for the use of LLMs, including recommendations for prompt engineering, evaluating the fairness and robustness of models, and more. Lastly, in order to evaluate the comprehensiveness and ease-of-use of the new pipeline, two replication studies will be conducted using the pipeline to test its ability to strengthen or contradict the findings of the original papers.

MSc. Artificial Intelligence
30EC

The Ethics and Privacy Quick Scan of the Utrecht University Research Institute of Information and Computing Sciences classified this research as low-risk with no fuller ethics review or privacy assessment required.

Contents

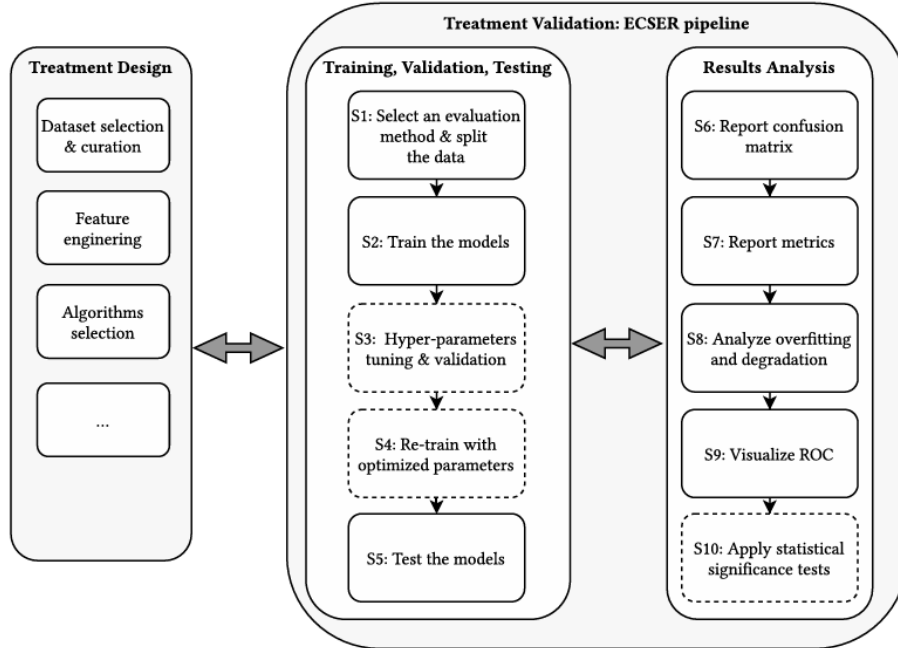
1	Introduction	3
2	Background	4
3	Research Approach	5
4	Related Work	8
4.1	Challenges in Empirical Software Engineering	8
4.2	Replications in SE	9
4.3	Large Language Models for SE	10
4.4	Prompt Engineering	11
4.5	Evaluating LLM Classifiers	11
	References	13

1 Introduction

In recent years, large language models (LLMs) have become an increasingly relevant tool in software engineering (SE), allowing researchers to accomplish complex tasks usually reserved for humans or other machine learning (ML) methods, such as code generation, summarization and classification. However, few guidelines exist for conducting and reporting on LLM research, leaving researchers on their own to decide what choices to make when designing prompts, what evaluation metrics to report in their results and the level of statistical validation needed for their comparisons. Because of these lack of guidelines, researchers using LLMs often omit important details, such as the F1 score or the significance of their results [1]. This problem is not new or unique to LLMs; many researchers have reported that there is a poor standard of empirical software engineering, which inhibits the usefulness of results and impairs reproducibility [2, 3].

To treat the observed poor standard of empirical SE, some researchers have called for a more systematic approach to conducting and evaluating experiments. The ECSER (Evaluating Classifiers in Software Engineering) pipeline, introduced by Dell’Anna et al. [4], provides a systemic pipeline for training, validating and testing ML classifiers and analysing their results. ECSER is based on recommendations from the field of machine learning for software engineering (ML4SE) and is designed to be easy to use for SE researchers, regardless of expertise with machine learning methods. The ECSER pipeline has seen some empirical validation in the form of two replication studies, in requirements engineering and software testing, where Dell’Anna et al. were able to confirm and strengthen some findings, as well as discover additional findings or findings that contradict the original ones. Figure 1 shows the steps of the ECSER pipeline.

Figure 1 The ECSER pipeline. Reproduced from Dell’Anna et al. [4]. Steps with a dashed border are optional. The \leftrightarrow arrows indicate feedback loops between phases



However, while ECSER is tailored to traditional ML classifiers (such as logistic regression, SVMs and decision trees), it includes no specific recommendations for LLM-based classifiers, despite LLMs showing

promise for various classification tasks [5, 6]. Thus, the ECSER pipeline is not sufficient for LLM research, since using LLMs for classification comes with unique considerations that are not present for traditional classifiers. For one, it is common to use pre-trained, generalized LLMs and give task-specific instructions to get the desired classification on the input data. This means that designing these instructions, also called prompts, becomes a crucial part to conducting LLM research [7]. There are also unique considerations when it comes to evaluating the fairness and robustness of the results: does the model behave in accordance with human values and are the results resilient to adversarial prompts, such as those containing typos [8, 9]?

The goal of this thesis is to address the lack of LLM-specific guidelines in the ECSER pipeline, by expanding it to include specific recommendations for conducting LLM classifier research. In doing so, we will make the following contributions:

- We will conduct a mapping study to assess the current state of reporting in LLM4SE classification research.
- We will develop an expanded version of the ECSER pipeline that includes specific recommendations for conducting LLM classifier research, such as prompt engineering and the evaluation of model fairness and robustness.
- We will conduct two replication studies of existing LLM classifier papers using the enhanced pipeline, in order to evaluate the pipeline’s usability and comprehensiveness.

The remainder of this thesis is structured as follows. In Section 2, we provide background information about LLMs. Section 3 describes the research approach and research questions. Lastly, Section 4 discusses related work.

2 Background

Language Models (LMs) are computational models designed to model the generative likelihood of word sequences [10]. In other words, given a sequence of input text, they predict the next word in the sequence based on what is deemed most likely by the model. One of the oldest examples of language models are N-grams, which model the likelihood of a word given the previous N tokens [11]. Early statistical models like N-gram models suffer from several problems, including limited context and difficulty with unseen words.

Pre-trained Language Models (PLMs) are LMs that are designed to take in vast amounts of unannotated text in order to learn knowledge about language and the world [12]. The knowledge that is contained in these models can then be used on new tasks by way of **transfer learning** [13]. PLMs require much more information to be contained within than statistical language models such as N-grams could provide, which has led researchers to design more sophisticated models that rely on neural networks to model language, starting with recurrent neural networks (RNNs) and later Long Short-Term Memory (LSTM) models [14, 15]. The biggest recent breakthrough came with the invention of **transformers**, which performed better than previous models, were easier to train and allowed for the parallel processing of inputs [16]. Transformers typically consist of an **encoder**, which embeds the input sequence into a hidden (latent) space and a **decoder**, which translates the abstract representation of the hidden space to the target text. Both the encoder and decoder use a self-attention mechanism to weigh the contribution of each token to the output [16, 17].

Large Language Models (LLMs) were introduced as a way to refer to PLMs with massive parameter sizes, since it was found that larger parameter sizes lead to better performance and emergent abilities [18, 19]. Almost all LLMs use transformer models with the self-attention mechanism, but not all of them use both the encoder and decoder component [20]:

- **Encoder-decoder LLMs**, such as BART [21] and T5X [22], are good at both understanding and generating language [20]. As such they are useful for tasks like summarization and translation [23, 24].
- **Encoder-only LLMs**, such as BERT and its variants [25–28], are tailored towards language understanding [20]. As such they are good at tasks like classification and inference [29].
- **Decoder-only LLMs**, such as the GPT-series [30–32], are tailored towards language generation [20]. As such they are good at all generation tasks, including code generation [33] and creating conversational agents like ChatGPT [34] and LLama 2 [35].

Prompt Engineering is the process of systematically designing natural language instructions (or **prompts**) to guide LLMs towards a desired output [36]. Prompts condition the outputs of the LLM on the given instructions, so that the model generates the next tokens with both the prompt and previously generated tokens in mind [12]. The most common approaches to prompt engineering are:

- **Zero-shot Prompting.** The model is given no examples for a task, forcing it to rely entirely on the prompt [37].
- **Few-shot Prompting.** A limited number of examples for a task are given, which the model can learn from [38].
- **Chain-of-Thought Prompting.** The prompt instructs the model to follow coherent reasoning steps [39].

Another less common approach is automatic prompt engineering, where instructions are automatically generated and selected [40].

3 Research Approach

The usage of LLMs in SE research has increased massively over the past few years, but the research of how to best conduct studies and report results in this field remains very limited. In order to approach the problem of LLM4SE reporting standards and design meaningful guidelines, we follow Wieringa’s design cycle of *problem investigation*, *treatment design* and *treatment validation* [41]. This leads us to the following research questions:

RQ1: What is the current state of reporting in LLM4SE classification research?

The purpose of RQ1 is both descriptive and evaluative: what do researchers report in their studies and does this level of reporting contribute or detract from the accuracy and reproducibility of the results. To answer the question, we will conduct a mapping study of SE papers from 2024 that use LLM classifiers using the guidelines by Kitchenham et al. [42, 43], summarizing what information researchers do (or do not) include in their papers.

Search Strategy

We will perform a manual search over three of the top SE conferences and journals (table 1). These venues were narrowed down from the list of venues included in [20]. The title and abstract of each paper will be checked to identify papers that are related to LLMs and classification, after which the full text will be checked to identify if the paper follows the inclusion and exclusion criteria (table 2).

Table 1 Publication venues included in the mapping study.

Acronym	Venue
ICSE	International Conference on Software Engineering
FSE	International Conference on the Foundations of Software Engineering
TSE	Transactions on Software Engineering

Table 2 Inclusion/Exclusion Criteria

Inclusion Criteria
The paper uses an LLM. The LLM is used for an SE classification task. Classification is done using prompts. The full text of the paper is accessible. It was published in 2024.
Exclusion Criteria
The paper fine-tunes an LLM. The paper does not use prompts for classification. The paper is a reprint or different version of another paper.

The information that will be extracted from the papers will be determined based on what is recommended to be included by the ECSE pipeline [4], metrics that are recommended by other authors (such as the Matthews correlation coefficient [44–46]) and metrics and other information relevant to LLMs, like prompting technique and measures of fairness & robustness (see section 4.4 and 4.5).

RQ2: What resource can be developed to assist researchers in conducting and reporting on LLM4SE classification research?

RQ2 takes the role of treatment design; how can we solve potential problems that are identified by RQ1? The purpose of RQ2 is to design an artifact that is easy to use and comprehensive. To answer the question, we will enhance the existing ECSE pipeline [4] to include the use of LLMs, taking into consideration recommendations from the fields of LLM4SE and ML4SE, as well as relevant statistical or methodological research in other fields (see section 4). Table 3 includes a preliminary summary of the additions that will be made to the ECSE pipeline, with the primary supporting references for the additions in brackets. The list of additions is subject to change based on the findings from RQ1 and further literature review. It is also possible that the best course of action will be to add more steps to the ECSE pipeline, as opposed to expanding existing ones, particularly for reporting the calibration, fairness and robustness or for prompt engineering.

Table 3 Preliminary summary of additions to the ECSER pipeline

ECSER Step	Additions
S1. Select an evaluation method and split the data	Analyse desired level of model alignment and safety [1]
S2. Train the model	Define the goal of the prompt [7, 47] Select prompt engineering approach (zero-shot, few-shot, etc.) [20, 47] Select prompting pattern [47, 48]
S3. Hyper-parameters tuning and validation	Tune LLM model temperature [49] Compare prompting patterns [36]
S7. Report Metrics	Report Matthews correlation coefficient (MCC) [44–46] Report calibration (e.g., expected calibration error [50]) Report fairness (e.g., equalized odds difference [8]) Report robustness (e.g., performance drop rate [9])

RQ3: How useful is the proposed pipeline in assisting LLM4SE classification research?

The purpose of RQ3 is treatment validation: does the designed pipeline in RQ2 achieve its goals of ease of use and comprehensiveness. To answer this question we will conduct two replication studies in different subfields of SE that employ LLM classification using the guidelines by Carver [51]. The original studies will be chosen based on the following criteria:

1. The authors provide a ready-to-use replication package.
2. The study is recent (past five years) and represents current practices in LLM4SE well.
3. The study is peer-reviewed and published in a reputable journal.
4. The study does not already follow all the proposed guidelines, so that the added benefit of the pipeline can be evaluated.

By conducting these replication studies, we can evaluate the comprehensiveness of the pipeline and identify potential shortcomings. In order to test the ease of use for the intended users, we aim to ask an SE researcher to use the pipeline or give their expert opinion on it if possible. After the first replication study, the findings from RQ3 will be taken into account to make changes to the guidelines from RQ2 where needed, creating a feedback-loop between RQ2 and RQ3. The second replication will be conducted using the updated pipeline and might inform further changes.

Timeline

The below table proposes a timeline for each step of the research project. It is subject to change if certain parts take more or less time than expected.

Table 4 Timeline

Period	Activity
Week 21–24	RQ1: Conduct and analyse mapping study
Week 25–26	RQ2: Design enhanced ECSER pipeline
Week 27–35	Summer break
Week 36–39	RQ2: Design enhanced ECSER pipeline (contd.)
Week 40–45	RQ3: Conduct and analyse replication studies
Week 46–48	Finish thesis
Week 49–50	Prepare and hold thesis defence

4 Related Work

The main focus of this thesis is to study the usage and evaluation of LLM Classifiers in Software Engineering, but in doing so we touch on several related research fields, including Empirical Software Engineering, Machine Learning for Software Engineering (ML4SE), Large Language Models for Software Engineering (LLM4SE) and Prompt Engineering. Additionally, each of the three research questions has related work that is relevant to their methodologies. Hence, we have split the related work into several sections based on different areas of study to aid readability.

First, we discuss challenges in empirical software engineering which motivated the creation of the ECSEER pipeline and further motivate our mapping study of LLM4SE. Second, we discuss the topic of replication in SE, which is relevant to our replication studies and the field of empirical SE. Third, we discuss how Large Language Models have been used in SE, with a focus on classification tasks. Fourth, we discuss related work in Prompt Engineering, as it has become a vital part to using LLMs and will inform our guidelines for prompt engineering in the enhanced ECSEER pipeline. Finally, we discuss related work on the evaluation of LLMs and classifiers which will inform our guidelines for the evaluation of LLM Classifiers.

4.1 Challenges in Empirical Software Engineering

Over the past couple decades, the use of machine learning methods has become ubiquitous in every field of science, among them software engineering. Zhang and Tsai [52] summarized seven main activities that use machine learning in SE: prediction, property/model discovery, transformation, generation, library construction and maintenance, acquisition of specifications and development knowledge management. In addition, they list specific tasks belonging to each activity, such as *predicting* software quality and *generating* test cases. The field of classification research falls almost exclusively in the prediction category.

With the increased popularity of machine learning methods, many software packages have been developed to allow anyone to use these methods without having to write any code or even understand how these methods work. This leads to many researchers using machine learning without fully understanding the underlying assumptions of the methods they are using and assuming the results are correct without proper (statistical) analysis. This problem is not new: Kitchenham et al. [2] reported at the start of this century that there was a poor standard of empirical software engineering, mostly caused by a lack of statistical expertise from both the authors and reviewers. They go on to provide extensive guidelines for designing and conducting experiments as well as analysing, presenting and interpreting results. In another paper, Kitchenham et al. [53] discuss the benefits of evidence-based software engineering approach by drawing an analogy to evidence-based medicine, but note that the infrastructure needed for widespread adoption of evidence-based SE isn't there yet, which means that SE experiments are vulnerable to subject and experimenter bias.

Methodological problems in SE research can result in problems with the reproducibility of results, which severely reduces the real-world usefulness of said results. Menzies and Shepperd [3] argue that the main goal of science is *conclusion stability*, which means that when you discover an effect, it holds true irrespective of the situation. However, they note that in the software engineering literature there is often as much evidence *for* a given effect, as *against*. This lack of conclusion stability is caused by sources of bias and variance that are introduced at various steps: sampling, pre-processing, training, algorithm selection and the analysis of results. If these sources of bias and variance are accurately analysed and reported, that would make it easier for the authors and other researchers to interpret the conclusions.

Despite being a known issue, the problems of improper analysis and reporting have persisted over the years. In a mapping study of the proceedings of the International Conference of Software Engineering (ICSE) from 2019 to 2021, Dell'Anna et al. [4] analysed 60 SE papers related to classification and found that most papers

(36) reported on the precision and recall of their experiment, fewer reported the F-score (27) and accuracy (24), but only 14 papers explicitly justified why they chose the reported metrics. The full confusion matrix was only included in six papers. Other metrics like the receiver operator characteristic (ROC) curve and the area under the curve (AUC) were mostly absent with 3 and 9 respective mentions. Of course, as we know from statistics, it is not enough to report a metric and claim it's better than that of other models, since it could be coincidental, so statistical significance testing is required to make such claims. However, only six papers analysed the statistical significance of their results. Evidently, the situation is dire, which is why Dell'Anna et al. developed the ECSER pipeline to give SE researchers a systematic approach to conducting classification experiments and reporting and analysing the results.

4.2 Replications in SE

Replication studies are a vital part of empirical software engineering, allowing researchers to investigate the effects of alternative values for important attributes, vary the strategy with which hypotheses are investigated and make up for certain threats to validity [54]. The importance of replication is not unique to software engineering, after all the so-called "replication crisis" was a massive topic of debate in psychology research in recent years, since it was found that many previous results in the field could not be replicated, even if the original findings were statistically significant [55–57]. Cruz et al. [58] conducted a systematic literature review of replication studies in SE between 2013 and 2018 and found that the number of published replication studies has been steadily increasing, though there are research gaps in certain fields, suggesting that the SE community is taking an increased interest in replicating previous studies.

Shull et al. [59] identify two types of replication studies: *exact replications*, in which the procedures of the original study are strictly followed; and *conceptual replications*, in which the research problem of the original study is evaluated using a different procedure. Additionally, exact replications are subdivided into *dependent* and *independent* replications. Dependent replications keep most of the conditions of the experiment (close to) the same, which allows one to test the effect of specific variables on the results; whereas independent replications vary the original conditions of the experiment significantly, which allows one to test the robustness of the original results under different conditions. Juristo and Vegas [60] underline the value of non-exact replications, which they define similarly to how Shull et al. define independent exact replications, by conducting a multiple-case replication study using a newly proposed four-phase process, to show that non-identical replications can be used to learn new information about the original results.

While there is a consensus in SE that replication is important, Shepperd et al. [61] show that replication studies often suffer from the same problems as original studies, such as incomplete reporting, low power or potential researcher bias. Notably, many replication studies in their survey did not report any details on the dispersion (e.g., variance) of the response variables. In another paper, Shepperd [62] shows that because of wide prediction intervals, almost all replications are confirmatory, which means the added knowledge is negligible. Because of these problems, they suggest that researchers should focus more on broader meta-analyses instead of replication studies. Other common problems in conducting replication research include difficulty in getting them published, lack of guidelines and the unavailability of replication packages [58, 63].

Communication can also play a factor in successful replication: Vegas et al. [64] investigated the role of communication in successful replication of previous experiments and found that there should be some level of communication, orally or in writing, between the previous authors and authors of the replication study, to avoid unnecessary changes to the experiments.

Given these challenges, guidelines for conducting replications are clearly needed. Despite this, the number of existing guidelines remains low. Guidelines were proposed by Carver [51], which are still commonly used in replications, but these seem to be the only replication guidelines for software engineering and were intended

as a starting point for future guidelines. Other artifacts do exist to aid replications: Gómez et al. [65] proposed a classification of replications in SE, differentiating between literal, operational and conceptual replications, that can be used to identify which changes can be made in each type of replication and understand the level or verification needed for that type.

Abualhaija et al. [66] proposed an artifact, referred to as the ID-Card, for summarizing papers in RE research that use NLP techniques. The intention behind the ID-Card is that all the relevant information needed for replication is presented in an easy to read format and can be used for both replication and educational purposes. Lastly, Brandt et al. [67] developed the "Replication Recipe" for psychology research, which lists 36 questions that should be addressed when conducting a replication, most of which are also directly applicable to empirical SE.

4.3 Large Language Models for SE

A recent development in SE research is the increased relevance and use of LLMs. Fan et al. [68] looked at all preprints on arXiv categorised under computer science whose title included "LLM", "Large Language Model" or "GPT" and found the number of SE papers that mentioned LLMs grew exponentially from 0 in 2019, 5 in 2020, to 181 in 2023. Naturally, this is only an approximation of the total number of preprints on arXiv that use LLMs, but it undoubtedly signifies a substantial increase. Fan et al. also conducted an extensive survey of how LLMs have been applied to various software development activities (such as code generation and testing) and research domains (such as human-computer interaction and education), but the survey is almost exclusively focused on the generative use of LLMs and does not go into depth on the use of LLMs in classification research. Many other good surveys and analyses of the use of LLMs in Software Engineering have been written, but most of them similarly do not mention any applications to classification tasks and instead exclusively focus on generation [69–71].

Surveys have also been written about specific domains of software engineering: Wang et al. [72] did a comprehensive review of 102 studies that use LLMs for software testing and found that LLMs have commonly been applied to various tasks, including test case generation, test input generation, debugging, program repair and more. They also found that while most papers used LLMs to address the entire task, many others combined LLMs with additional techniques to optimize the outputs of the LLM, including mutation testing, differential testing, syntactic checking, program analysis, statistical analysis and other techniques. With these extra techniques, researchers were able to generate more diverse and complex code and overcome some of the limitations of LLMs.

While LLMs are most frequently used to generate text (and code), they have seen considerable success in classification tasks. Guo et al. [5] have shown that LLMs can outperform other methods like SVMs for health-related text classification. Fields et al. [6] present an in-depth survey of text classification using transformers across domains and found that LLMs can perform remarkably well on many (but not all) classification tasks. However, Chen et al. [73] compared the performance of various LLMs to traditional ML methods in clinical prediction and found that LLMs could *not* beat traditional methods in this case. Vajjala and Shimangaud [74] also show that there are large performance disparities between languages in classification tasks. Thus, LLMs might not always be a better choice than traditional ML methods, but they certainly show potential.

Software engineers have also started using LLMs for classification. Hou et al. [20] conducted a systematic literature review of 395 software engineering studies that use LLMs and found that around 21.61% involve classification tasks. The most common classification tasks where LLMs have been used include vulnerability detection, requirements classification, bug prediction and review/commit/code classification, with many other classification tasks having been attempted using LLMs. Zhang et al. [75] conducted a systematic survey of 947 SE studies and summarized 62 unique LLMs of code, including six LLMs specifically fine-tuned for the

tasks of code classification, clone detection and defect detection.

4.4 Prompt Engineering

Designing a good prompt for a particular task is vital to ensuring the adequate performance of LLMs, which is why much research has gone into the systematic designing of prompts, or prompt engineering. Marvin et al. [7] provide an overview of prompt engineering principles and techniques and outline the main steps involved in the process of prompt engineering: 1) define the goal of the prompt; 2) understand the model’s capabilities; 3) choose the right prompting format; 4) provide context to the LLM and 5) test and refine the prompt based on the goal. Several resources exist to make it easier to choose the right prompting approach: White et al. [48] introduce prompt patterns, which offer reusable solutions to common prompting challenges, and provide a framework for designing and documenting these patterns in terms of the intent, motivation, structure and consequences of the patterns. For example, the persona pattern can be used to make the LLM take the point of view of an expert in the field, which is especially useful if the user itself is not an expert. The paper provides a catalogue of this pattern and other successfully applied patterns with example implementations. Further, Ekin [76] provides an accessible (AI-generated) guide to prompt engineering.

Prompt engineering research has also been conducted specifically for software engineering: Hou et al. [20] looked at what prompt engineering techniques are commonly applied in SE tasks and found the most common techniques involve few-shot prompting and zero-shot prompting, but the third largest group of studies had no explicit mention of prompting techniques or proposed their own strategies. Other techniques included chain-of-thought, automatic prompt engineering, chain-of-code, automatic chain-of-thought, modular-of-thought and structured chain-of-thought.

Arvidsson and Axell [47] collected prompt engineering recommendations for requirements engineering (RE) in a systematic literature review and classified the guidelines into various themes, after which they interviewed three RE experts to evaluate the guidelines. Ronanki et al. [36] evaluated the effectiveness of 5 prompting patterns (from the catalogue by White et al. [48]) on two RE tasks and proposed a framework for evaluating the effectiveness of prompting patterns for any RE task, providing five steps for conducting a comparison of patterns.

A recent development in prompt engineering is automatic prompt engineering (APE), which has shown good performance compared to baseline models, while avoiding the need for manually designing prompts [40, 77]. Zadenoori et al. [78] investigated the use of automatic prompt engineering in RE and found that, on average, APE outperforms traditional prompt engineering techniques. However, research is still limited.

4.5 Evaluating LLM Classifiers

Few papers have been written specifically about the evaluation of LLM-based classifiers, but we can look at the evaluation of classifiers and LLMs separately to get a picture. The ECSER pipeline is the most relevant for evaluating classifiers, since it gives in-depth recommendations about the conducting and reporting of classifier research [4]. Specific recommendations include the reporting of the full confusion matrix, reporting other metrics relevant to the domain such as specificity (true negative rate), analysing overfitting and degradation, visualising the ROC and applying statistical significance tests. Some researchers also recommend reporting the Matthews correlation coefficient (MCC), which combines all four metrics of the confusion matrix and can be more informative than the F1 score [44–46].

Hou et al. [20] looked at whether these metrics were reported in LLM-based classification research and found that out of 147 papers that use LLMs for classification, the most frequently reported metric is precision with 35 papers, followed by recall (34), F1-score (33) and accuracy (23), but the AUC was reported 9 times,

the ROC only 4 times and the MCC only twice. Unfortunately, they did not count how many studies did significance testing. These results are very similar to Dell’Anna et al. [4]’s aforementioned mapping study of classifier research as a whole and suggest that many important metrics are under-reported in LLM-based classification research.

Using and evaluation LLMs comes with more considerations than just the accuracy of the results. Guo et al. [1] categorize the evaluation of LLMs into three types: knowledge and capability evaluation, which assesses the fundamental knowledge and reasoning capabilities of the LLM; alignment evaluation, which refers to evaluating ethical and moral considerations like bias; and safety evaluation, which focuses on the robustness of LLMs and risk evaluation. Liu et al. [79] have created guidelines specifically for evaluating LLM alignment in terms of reliability, safety, fairness, resistance to misuse, explainability and reasoning, adherence to social norms and robustness.

Another consideration is how well the LLM outputs adhere to the format specified in the prompt. For example, it is common to ask an LLM to format its answer as a valid JSON object, allowing for easier processing of the output. Limiting responses to a specified format can be beneficial: Tam et al. [80] compared the performance of several LLMs between free-form response and formatted responses and found that, while reasoning ability was weakened for formatted responses, classification accuracy was increased. However, LLMs do not always follow the format accurately. Long et al. [81] define several measures to analyse the level of adherence to the output format: SysE, which measures the performance of the LLM for answers that *meet the constraint*; TrueE, which measures the performance of all answers *regardless of whether the format is satisfied*; and BiasF, which measures the (mean squared) difference between SysE and TrueE. Li et al. [82] introduce the JScore measure, which measures the similarity between JSON objects and can be used to compare LLM generated JSON objects with the target JSON objects. Lastly, Xia et al. [83] introduce the FoFo benchmark for evaluating LLMs based on their ability to follow complex, domain-specific formats.

Chang et al. [84] conducted another survey of the evaluation of LLMs which focuses on downstream applications, resulting in a classification of four aspects of LLM evaluation: the accuracy, containing measures of correctness such as the F1 score; the calibration, which contains measures for the degree of alignment between the predicted probabilities and actual probabilities, such as the expected calibration error (ECE); the fairness, with measures pertaining to the equal treatment of different groups, such as the equalized odds difference (EOD); and robustness, with measures pertaining to the performance of a model in the face of challenging inputs and noise, such as performance drop rate (PDR). Further, they mention several ways in which human evaluation can be used to evaluate LLMs, such as the degree to which the model outputs align with human values.

Lastly, the model temperature, which regulates the randomness of the model, can also affect its performance. Peeperkorn et al. [49] looked at the effect of model temperature on the level of creativity of the outputs and found that LLMs generate slightly more novel outputs with higher temperature, but the temperature was also correlated with incoherence. No relationship was found between temperature and cohesion or typicality. Model temperature might also affect a model’s susceptibility to security attacks: Yu et al. [85] found that some LLMs became more susceptible to jailbreaking attacks as temperature increased, whereas others had decreased susceptibility to jailbreaking with increased temperature, possibly due to the fact that the former models had a lower susceptibility at 0 temperature and the latter models had a higher susceptibility at 0 temperature.

References

- [1] Z. Guo et al., *Evaluating large language models: A comprehensive survey*, 2023. arXiv: 2310.19736 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/2310.19736>
- [2] B. Kitchenham et al., “Preliminary guidelines for empirical research in software engineering,” *IEEE Transactions on Software Engineering*, vol. 28, no. 8, pp. 721–734, 2002. DOI: 10.1109/TSE.2002.1027796
- [3] T. Menzies and M. Shepperd, “Special issue on repeatable results in software engineering prediction,” *Empirical Software Engineering*, vol. 17, no. 1–2, pp. 1–17, Jan. 2012, ISSN: 1573-7616. DOI: 10.1007/s10664-011-9193-5 [Online]. Available: <http://dx.doi.org/10.1007/s10664-011-9193-5>
- [4] D. Dell’Anna, F. B. Aydemir, and F. Dalpiaz, “Evaluating classifiers in se research: The ecser pipeline and two replication studies,” *Empirical Softw. Engg.*, vol. 28, no. 1, Nov. 2022, ISSN: 1382-3256. DOI: 10.1007/s10664-022-10243-1 [Online]. Available: <https://doi.org/10.1007/s10664-022-10243-1>
- [5] Y. Guo, A. Ovadje, M. Al-garadi, and A. Sarker, “Evaluating large language models for health-related text classification tasks with public social media data,” *Journal of the American Medical Informatics Association : JAMIA*, vol. 31, pp. 2181–2189, 2024. DOI: 10.1093/jamia/ocae210
- [6] J. Fields, K. Chovanec, and P. Madiraju, “A survey of text classification with transformers: How wide? how large? how long? how accurate? how expensive? how safe?” *IEEE Access*, vol. 12, pp. 6518–6531, 2024. DOI: 10.1109/ACCESS.2024.3349952
- [7] G. Marvin, N. Hellen, D. Jjingo, and J. Nakatumba-Nabende, “Prompt engineering in large language models,” in *Data Intelligence and Cognitive Informatics*, I. J. Jacob, S. Piramuthu, and P. Falkowski-Gilski, Eds., Singapore: Springer Nature Singapore, 2024, pp. 387–402, ISBN: 978-981-99-7962-2.
- [8] B. Woodworth, S. Gunasekar, M. I. Ohannessian, and N. Srebro, “Learning non-discriminatory predictors,” in *Proceedings of the 2017 Conference on Learning Theory*, S. Kale and O. Shamir, Eds., ser. Proceedings of Machine Learning Research, vol. 65, PMLR, Jul. 2017, pp. 1920–1953. [Online]. Available: <https://proceedings.mlr.press/v65/woodworth17a.html>
- [9] K. Zhu et al., *Promptrobust: Towards evaluating the robustness of large language models on adversarial prompts*, 2024. arXiv: 2306.04528 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/2306.04528>
- [10] W. X. Zhao et al., *A survey of large language models*, 2025. arXiv: 2303.18223 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/2303.18223>
- [11] C. E. Shannon, “A mathematical theory of communication,” *The Bell System Technical Journal*, vol. 27, no. 3, pp. 379–423, 1948. DOI: 10.1002/j.1538-7305.1948.tb01338.x
- [12] D. Jurafsky and J. H. Martin, “Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition with language models,” Online manuscript released January 12, 2025, 2025. [Online]. Available: <https://web.stanford.edu/~jurafsky/slp3/>
- [13] H. Wang, J. Li, H. Wu, E. Hovy, and Y. Sun, “Pre-trained language models and their applications,” *Engineering*, vol. 25, pp. 51–65, 2023, ISSN: 2095-8099. DOI: <https://doi.org/10.1016/j.eng.2022.04.024> [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S2095809922006324>
- [14] T. Mikolov and G. Zweig, “Context dependent recurrent neural network language model,” in *2012 IEEE Spoken Language Technology Workshop (SLT)*, IEEE, 2012, pp. 234–239.
- [15] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

- [16] A. Vaswani et al., *Attention is all you need*, 2017. DOI: 10.48550/ARXIV.1706.03762 [Online]. Available: <https://arxiv.org/abs/1706.03762>
- [17] D. Bahdanau, K. Cho, and Y. Bengio, *Neural machine translation by jointly learning to align and translate*, 2016. arXiv: 1409.0473 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/1409.0473>
- [18] J. Kaplan et al., *Scaling laws for neural language models*, 2020. arXiv: 2001.08361 [cs.LG]. [Online]. Available: <https://arxiv.org/abs/2001.08361>
- [19] J. Wei et al., *Emergent abilities of large language models*, 2022. arXiv: 2206.07682 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/2206.07682>
- [20] X. Hou et al., *Large language models for software engineering: A systematic literature review*, 2024. arXiv: 2308.10620 [cs.SE]. [Online]. Available: <https://arxiv.org/abs/2308.10620>
- [21] M. Lewis et al., *Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension*, 2019. arXiv: 1910.13461 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/1910.13461>
- [22] A. Roberts et al., “Scaling up models and data with t5x and seqio,” *arXiv preprint arXiv:2203.17189*, 2022. [Online]. Available: <https://arxiv.org/abs/2203.17189>
- [23] K. Cho, B. van Merriënboer, D. Bahdanau, and Y. Bengio, *On the properties of neural machine translation: Encoder-decoder approaches*, 2014. arXiv: 1409.1259 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/1409.1259>
- [24] A. Asadi and R. Safabakhsh, “The encoder-decoder framework and its applications,” in *Deep Learning: Concepts and Architectures*, W. Pedrycz and S.-M. Chen, Eds. Cham: Springer International Publishing, 2020, pp. 133–167, ISBN: 978-3-030-31756-0. DOI: 10.1007/978-3-030-31756-0_5 [Online]. Available: https://doi.org/10.1007/978-3-030-31756-0_5
- [25] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, *Bert: Pre-training of deep bidirectional transformers for language understanding*, 2019. arXiv: 1810.04805 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/1810.04805>
- [26] Y. Liu et al., *Roberta: A robustly optimized bert pretraining approach*, 2019. arXiv: 1907.11692 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/1907.11692>
- [27] V. Sanh, L. Debut, J. Chaumond, and T. Wolf, “Distilbert, a distilled version of bert: Smaller, faster, cheaper and lighter,” *arXiv preprint arXiv:1910.01108*, 2019.
- [28] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut, *Albert: A lite bert for self-supervised learning of language representations*, 2020. arXiv: 1909.11942 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/1909.11942>
- [29] M. V. Koroteyev, *Bert: A review of applications in natural language processing and understanding*, 2021. arXiv: 2103.11943 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/2103.11943>
- [30] A. Radford, K. Narasimhan, T. Salimans, I. Sutskever, et al., “Improving language understanding by generative pre-training,” 2018.
- [31] T. B. Brown et al., *Language models are few-shot learners*, 2020. arXiv: 2005.14165 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/2005.14165>
- [32] OpenAI et al., *Gpt-4 technical report*, 2024. arXiv: 2303.08774 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/2303.08774>
- [33] R. A. Poldrack, T. Lu, and G. Beguš, *Ai-assisted coding: Experiments with gpt-4*, 2023. arXiv: 2304.13187 [cs.AI]. [Online]. Available: <https://arxiv.org/abs/2304.13187>
- [34] OpenAI, *Chatgpt*, <https://chat.openai.com/>, Available at <https://chat.openai.com/>, 2023.
- [35] H. Touvron et al., *Llama 2: Open foundation and fine-tuned chat models*, 2023. arXiv: 2307.09288 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/2307.09288>

- [36] K. Ronanki, B. Cabrero-Daniel, J. Horkoff, and C. Berger, *Requirements engineering using generative ai: Prompts and prompting patterns*, 2023. arXiv: 2311.03832 [cs.SE]. [Online]. Available: <https://arxiv.org/abs/2311.03832>
- [37] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, et al., “Language models are unsupervised multitask learners,” *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.
- [38] T. Brown et al., “Language models are few-shot learners,” *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.
- [39] J. Wei et al., “Chain-of-thought prompting elicits reasoning in large language models,” *Advances in neural information processing systems*, vol. 35, pp. 24 824–24 837, 2022.
- [40] Y. Zhou et al., *Large language models are human-level prompt engineers*, 2023. arXiv: 2211.01910 [cs.LG]. [Online]. Available: <https://arxiv.org/abs/2211.01910>
- [41] R. J. Wieringa, *Design Science Methodology for Information Systems and Software Engineering*. Springer Berlin Heidelberg, 2014, ISBN: 9783662438398. DOI: 10.1007/978-3-662-43839-8 [Online]. Available: <http://dx.doi.org/10.1007/978-3-662-43839-8>
- [42] B. A. Kitchenham, D. Budgen, and O. Pearl Brereton, “Using mapping studies as the basis for further research – a participant-observer case study,” *Information and Software Technology*, vol. 53, no. 6, pp. 638–651, 2011, Special Section: Best papers from the APSEC, ISSN: 0950-5849. DOI: <https://doi.org/10.1016/j.infsof.2010.12.011> [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0950584910002272>
- [43] B. Kitchenham and S. Charters, “Guidelines for performing systematic literature reviews in software engineering,” vol. 2, Jan. 2007.
- [44] D. Chicco and G. Jurman, “The advantages of the matthews correlation coefficient (mcc) over f1 score and accuracy in binary classification evaluation,” *BMC Genomics*, vol. 21, no. 1, Jan. 2020, ISSN: 1471-2164. DOI: 10.1186/s12864-019-6413-7 [Online]. Available: <http://dx.doi.org/10.1186/s12864-019-6413-7>
- [45] G. M. Foody, “Challenges in the real world use of classification accuracy metrics: From recall and precision to the matthews correlation coefficient,” *PLOS ONE*, vol. 18, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:263668955>
- [46] J. Yao and M. J. Shepperd, “Assessing software defection prediction performance: Why using the matthews correlation coefficient matters,” *Proceedings of the 24th International Conference on Evaluation and Assessment in Software Engineering*, 2020. [Online]. Available: <https://api.semanticscholar.org/CorpusID:211818318>
- [47] S. Arvidsson and J. Axell, “Prompt engineering guidelines for llms in requirements engineering,” 2023.
- [48] J. White et al., *A prompt pattern catalog to enhance prompt engineering with chatgpt*, 2023. arXiv: 2302.11382 [cs.SE]. [Online]. Available: <https://arxiv.org/abs/2302.11382>
- [49] M. Peeperkorn, T. Kouwenhoven, D. Brown, and A. Jordanous, *Is temperature the creativity parameter of large language models?* 2024. arXiv: 2405.00492 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/2405.00492>
- [50] M. Pakdaman Naeini, G. Cooper, and M. Hauskrecht, “Obtaining well calibrated probabilities using bayesian binning,” *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 29, no. 1, Feb. 2015. DOI: 10.1609/aaai.v29i1.9602 [Online]. Available: <https://ojs.aaai.org/index.php/AAAI/article/view/9602>
- [51] J. C. Carver, “Towards reporting guidelines for experimental replications: A proposal,” in *Proceedings of the 1st International Workshop on Replication in Empirical Software Engineering Research (RESER)*, Held during ICSE 2010, Cape Town, South Africa, May 2010.
- [52] D. Zhang and J. Tsai, “Machine learning and software engineering,” vol. 11, Feb. 2002, pp. 22–29, ISBN: 0-7695-1849-4. DOI: 10.1109/TAI.2002.1180784

- [53] B. Kitchenham, T. Dybå, and M. Jorgensen, “Evidence-based software engineering,” Jun. 2004, pp. 273–281, ISBN: 0-7695-2163-0. DOI: 10.1109/ICSE.2004.1317449
- [54] V. Basili, F. Shull, and F. Lanubile, “Building knowledge through families of experiments,” *IEEE Transactions on Software Engineering*, vol. 25, no. 4, pp. 456–473, 1999. DOI: 10.1109/32.799939
- [55] P. E. Shrout and J. L. Rodgers, “Psychology, science, and knowledge construction: Broadening perspectives from the replication crisis,” *Annual Review of Psychology*, vol. 69, no. Volume 69, 2018, pp. 487–510, 2018, ISSN: 1545-2085. DOI: <https://doi.org/10.1146/annurev-psych-122216-011845> [Online]. Available: <https://www.annualreviews.org/content/journals/10.1146/annurev-psych-122216-011845>
- [56] S. E. Maxwell, M. Y.-K. Lau, and G. S. Howard, “Is psychology suffering from a replication crisis? what does ”failure to replicate” really mean?” *The American psychologist*, vol. 70 6, pp. 487–98, 2015. [Online]. Available: <https://api.semanticscholar.org/CorpusID:11690879>
- [57] V. Amrhein, D. Trafimow, and S. Greenland, “Inferential statistics as descriptive statistics: There is no replication crisis if we don’t expect replication,” *The American Statistician*, vol. 73, pp. 262–270, 2018. [Online]. Available: <https://api.semanticscholar.org/CorpusID:56130495>
- [58] M. Cruz, B. Bernárdez, A. Durán, J. A. Galindo, and A. Ruiz-Cortés, “Replication of studies in empirical software engineering: A systematic mapping study, from 2013 to 2018,” *IEEE Access*, vol. 8, pp. 26 773–26 791, 2020. [Online]. Available: <https://api.semanticscholar.org/CorpusID:209073520>
- [59] F. J. Shull, J. C. Carver, S. Vegas, and N. Juristo, “The role of replications in empirical software engineering,” *Empirical Software Engineering*, vol. 13, no. 2, pp. 211–218, Jan. 2008, ISSN: 1573-7616. DOI: 10.1007/s10664-008-9060-1 [Online]. Available: <http://dx.doi.org/10.1007/s10664-008-9060-1>
- [60] N. Juristo and S. Vegas, “The role of non-exact replications in software engineering experiments,” *Empirical Software Engineering*, vol. 16, no. 3, pp. 295–324, Aug. 2010, ISSN: 1573-7616. DOI: 10.1007/s10664-010-9141-9 [Online]. Available: <http://dx.doi.org/10.1007/s10664-010-9141-9>
- [61] M. Shepperd, N. Ajenka, and S. Counsell, “The role and value of replication in empirical software engineering results,” *Information and Software Technology*, vol. 99, pp. 120–132, 2018, ISSN: 0950-5849. DOI: <https://doi.org/10.1016/j.infsof.2018.01.006> [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0950584917304305>
- [62] M. Shepperd, “Replication studies considered harmful,” in *Proceedings of the 40th International Conference on Software Engineering: New Ideas and Emerging Results*, ser. ICSE-NIER ’18, Gothenburg, Sweden: Association for Computing Machinery, 2018, pp. 73–76, ISBN: 9781450356626. DOI: 10.1145/3183399.3183423 [Online]. Available: <https://doi.org/10.1145/3183399.3183423>
- [63] J. Siegmund, N. Siegmund, and S. Apel, “Views on internal and external validity in empirical software engineering,” in *2015 IEEE/ACM 37th IEEE International Conference on Software Engineering*, vol. 1, 2015, pp. 9–19. DOI: 10.1109/ICSE.2015.24
- [64] S. Vegas, N. Juristo, A. Moreno, M. Solari, and P. Letelier, “Analysis of the influence of communication between researchers on experiment replication,” in *Proceedings of the 2006 ACM/IEEE International Symposium on Empirical Software Engineering*, ser. ISESE ’06, Rio de Janeiro, Brazil: Association for Computing Machinery, 2006, pp. 28–37, ISBN: 1595932186. DOI: 10.1145/1159733.1159741 [Online]. Available: <https://doi.org/10.1145/1159733.1159741>
- [65] O. S. Gómez, N. Juristo, and S. Vegas, “Understanding replication of experiments in software engineering: A classification,” *Information and Software Technology*, vol. 56, no. 8, pp. 1033–1048, 2014, ISSN: 0950-5849. DOI: <https://doi.org/10.1016/j.infsof.2014.04.004>

- [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0950584914000858>
- [66] S. Abualhaija et al., “Replication in requirements engineering: The nlp for re case,” *ACM Trans. Softw. Eng. Methodol.*, vol. 33, no. 6, Jun. 2024, ISSN: 1049-331X. DOI: 10.1145/3658669 [Online]. Available: <https://doi.org/10.1145/3658669>
 - [67] M. J. Brandt et al., “The replication recipe: What makes for a convincing replication?” *Journal of Experimental Social Psychology*, vol. 50, pp. 217–224, 2014, ISSN: 0022-1031. DOI: <https://doi.org/10.1016/j.jesp.2013.10.005> [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0022103113001819>
 - [68] A. Fan et al., “Large language models for software engineering: Survey and open problems,” in *2023 IEEE/ACM International Conference on Software Engineering: Future of Software Engineering (ICSE-FoSE)*, 2023, pp. 31–53. DOI: 10.1109/ICSE-FoSE59343.2023.00008
 - [69] L. Belzner, T. Gabor, and M. Wirsing, “Large language model assisted software engineering: Prospects, challenges, and a case study,” in *Bridging the Gap Between AI and Reality*, B. Steffen, Ed., Cham: Springer Nature Switzerland, 2024, pp. 355–374, ISBN: 978-3-031-46002-9.
 - [70] Z. Zheng et al., “Towards an understanding of large language models in software engineering tasks,” *Empirical Software Engineering*, vol. 30, no. 2, Dec. 2024, ISSN: 1573-7616. DOI: 10.1007/s10664-024-10602-0 [Online]. Available: <http://dx.doi.org/10.1007/s10664-024-10602-0>
 - [71] I. Ozkaya, “Application of large language models to software engineering tasks: Opportunities, risks, and implications,” *IEEE Software*, vol. 40, no. 3, pp. 4–8, 2023. DOI: 10.1109/MS.2023.3248401
 - [72] J. Wang, Y. Huang, C. Chen, Z. Liu, S. Wang, and Q. Wang, *Software testing with large language models: Survey, landscape, and vision*, 2024. arXiv: 2307.07221 [cs.SE]. [Online]. Available: <https://arxiv.org/abs/2307.07221>
 - [73] C. Chen et al., “Clinicalbench: Can llms beat traditional ml models in clinical prediction?” *ArXiv*, vol. abs/2411.06469, 2024. [Online]. Available: <https://api.semanticscholar.org/CorpusID:273963111>
 - [74] S. Vajjala and S. Shimangaud, *Text classification in the llm era – where do we stand?* 2025. arXiv: 2502.11830 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/2502.11830>
 - [75] Q. Zhang et al., *A survey on large language models for software engineering*, 2024. arXiv: 2312.15223 [cs.SE]. [Online]. Available: <https://arxiv.org/abs/2312.15223>
 - [76] S. Ekin, “Prompt engineering for chatgpt: A quick guide to techniques, tips, and best practices,” 2023. DOI: 10.36227/techrxiv.22683919.v2
 - [77] Q. Ye, M. Axmed, R. Pryzant, and F. Khani, *Prompt engineering a prompt engineer*, 2024. arXiv: 2311.05661 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/2311.05661>
 - [78] M. A. Zadenoori, L. Zhao, W. Alhoshan, and A. Ferrari, “Automatic prompt engineering: The case of requirements classification,” in *International Working Conference on Requirements Engineering: Foundation for Software Quality*, Springer, 2025, pp. 217–225.
 - [79] Y. Liu et al., *Trustworthy llms: A survey and guideline for evaluating large language models’ alignment*, 2024. arXiv: 2308.05374 [cs.AI]. [Online]. Available: <https://arxiv.org/abs/2308.05374>
 - [80] Z. R. Tam, C.-K. Wu, Y.-L. Tsai, C.-Y. Lin, H.-y. Lee, and Y.-N. Chen, “Let me speak freely? a study on the impact of format restrictions on large language model performance.” in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*, F. Dernoncourt, D. Preoŕiuc-Pietro, and A. Shimorina, Eds., Miami, Florida, US: Association for Computational Linguistics, Nov. 2024, pp. 1218–1236. DOI: 10.18653/v1/2024.emnlp-industry.91 [Online]. Available: <https://aclanthology.org/2024.emnlp-industry.91/>

- [81] D. X. Long et al., *Llms are biased towards output formats! systematically evaluating and mitigating output format bias of llms*, 2025. arXiv: 2408.08656 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/2408.08656>
- [82] D. Li, Y. Zhao, Z. Wang, C. Jung, and Z. Zhang, "Large language model-driven structured output: A comprehensive benchmark and spatial data generation framework," *ISPRS International Journal of Geo-Information*, vol. 13, no. 11, 2024, ISSN: 2220-9964. DOI: 10.3390/ijgi13110405 [Online]. Available: <https://www.mdpi.com/2220-9964/13/11/405>
- [83] C. Xia et al., "FOFO: A benchmark to evaluate LLMs' format-following capability," in *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, L.-W. Ku, A. Martins, and V. Srikumar, Eds., Bangkok, Thailand: Association for Computational Linguistics, Aug. 2024, pp. 680–699. DOI: 10.18653/v1/2024.acl-long.40 [Online]. Available: <https://aclanthology.org/2024.acl-long.40/>
- [84] Y. Chang et al., *A survey on evaluation of large language models*, 2023. arXiv: 2307.03109 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/2307.03109>
- [85] C. X. Yu, J. James, C. Si Yu, D. David, C. B. Hoe, and P. Hui-Li Phyllis, "Can llms have a fever? investigating the effects of temperature on llm security," in *Proceedings of the International Conference on Industrial Engineering and Operations Management*, ser. South American '24, IEOM Society International, May 2024. DOI: 10.46254/sa05.20240024 [Online]. Available: <http://dx.doi.org/10.46254/SA05.20240024>