

Which Prompting Technique Should I Use? An Empirical Investigation of Prompting Techniques for Software Engineering Tasks

Enio Garcia Santana Junior¹, Gabriel Benjamin¹, Melissa Araujo¹, Harrison Santos¹, David Freitas¹, Eduardo Almeida¹, Paulo Anselmo da Mota Silveira Neto², Jiawei Li³, Jina Chun³, Iftekhar Ahmed³

¹Federal University of Bahia (UFBA), Brazil

²Federal Rural University of Pernambuco (UFRPE), Brazil

³University of California, Irvine (UCI), USA

Abstract—A growing variety of prompt engineering techniques has been proposed for Large Language Models (LLMs), yet systematic evaluation of each technique on individual software engineering (SE) tasks remains underexplored. In this study, we present a systematic evaluation of 14 established prompt techniques across 10 SE tasks using four LLM models. As identified in the prior literature, the selected prompting techniques span six core dimensions (Zero-Shot, Few-Shot, Thought Generation, Ensembling, Self-Criticism, and Decomposition). They are evaluated on tasks such as code generation, bug fixing, and code-oriented question answering, to name a few. Our results show which prompting techniques are most effective for SE tasks requiring complex logic and intensive reasoning versus those that rely more on contextual understanding and example-driven scenarios. We also analyze correlations between the linguistic characteristics of prompts and the factors that contribute to the effectiveness of prompting techniques in enhancing performance on SE tasks. Additionally, we report the time and token consumption for each prompting technique when applied to a specific task and model, offering guidance for practitioners in selecting the optimal prompting technique for their use cases.

Index Terms—Prompt Engineering; Large Language Models; Software Engineering Tasks.

I. INTRODUCTION

Emerging Large Language Models (LLMs) have rapidly advanced beyond the capabilities of earlier smaller models, thanks to breakthroughs in deep neural architectures, access to large-scale training corpora in natural language and source code, and substantial computational power [1]. Simply through prompting without any model training, LLMs have already shown state-of-the-art performance in automating various Software Engineering (SE) tasks, such as code translation [2], [3], commit message generation [4], [5], and program repair [6], [7]. Despite these promising developments, it has become increasingly evident that even small changes in how a prompt is formulated can drastically alter an LLM’s output [8]. In automating SE tasks by prompting LLMs, the adopted prompts have the potential to impact various quality aspects ranging from code correctness and readability to the efficiency of bug-fixing suggestions [9]–[11].

To guide LLMs toward more reliable, accurate, and context-aware responses, an array of *prompt engineering* [9] tech-

niques have been proposed by researchers in recent years. However, it remains unknown which prompt engineering technique would benefit a specific SE task the most, leaving practitioners uncertain about how to compose the prompts to suit their needs. In this work, we address this gap by *systematically evaluating the effectiveness of fourteen widely used prompt engineering techniques* [9] within *ten SE tasks* [12], spanning an extensive range of code understanding and generation tasks. This prompted us to ask our first research question:

RQ1: How do different prompting techniques impact the performance of SE tasks?

We analyzed the linguistic features of these prompts to gain a deeper understanding of the characteristics that distinguish prompting techniques that lead to performance improvements from those that do not in different SE tasks. This includes examining aspects such as Lexical Diversity, Token Count, Flesch-Kincaid Grade Level, Gunning Fog Index, and Flesch Reading Ease. These linguistic features help identify how prompt clarity and complexity relate to LLM performance in SE tasks. Lexical Diversity reflects vocabulary richness, Token Count captures prompt length, and Flesch-Kincaid Grade Level, Gunning Fog Index, and Flesch Reading Ease assess readability and complexity. Together, they reveal whether more concise, readable, or varied prompts are associated with improved task outcomes. Identifying significant linguistic similarities among high-performing prompting strategies can reveal useful insights for systematically designing more effective prompts for SE automation. In addition to linguistic features, we investigated other potential factors contributing to the success of prompting techniques by leveraging contrastive explanation [13]. Specifically, we queried the LLMs to compare successful and less successful prompts, aiming to understand what the model perceives as key differentiators in performance. Thus, we formulated the following research questions:

RQ2: What linguistic features of prompting techniques are associated with improved performance on SE tasks?

RQ3: What factors, according to LLMs, contribute to the effectiveness of prompting techniques in SE tasks?

Since large language models (LLMs) require substantial computational resources and can incur high financial costs, particularly when accessed via commercial APIs, software practitioners must carefully balance the trade-off between performance improvements and resource consumption [14]. Recent studies indicate that seemingly minor variations in prompts can result in disproportionately large differences in resource usage, including increased token counts and longer inference times. Therefore, understanding how different prompting techniques impact both performance and resource costs is essential for making informed decisions in practical SE scenarios. This motivation leads to our next research question:

RQ4: How are resource costs associated with the performance of prompting techniques in SE tasks?

Based on our findings, we constructed composite prompts by combining the prompting techniques that consistently yield performance improvements in SE tasks. We aimed to investigate whether integrating all positively impactful factors could lead to further performance gains. Ultimately, we aim to enable SE practitioners and researchers to leverage LLMs more effectively and accurately, advancing more robust and autonomous SE practices. This motivated our final research question:

Overall, this paper makes the following contributions:

- A comprehensive benchmark of **14 prompt-engineering techniques** evaluated on **10 SE tasks** and **3 LLM families**, totaling more than 2k prompts.
- A **systematic investigation** that discloses the factors associated with a prompting technique’s success across tasks and models.
- A **cost-aware analysis** reporting token budgets and latency for every technique, enabling practical cost–performance trade-offs.
- Public release of all datasets, prompts, model outputs, and analysis scripts to facilitate replication and future research.

II. RELATED WORK

LLMs in Software Engineering. Recent LLMs have billions of parameters and are trained with large-scale datasets of natural language and source code, showing state-of-the-art understanding and generation capabilities. Consequently, researchers have increasingly adopted them to conduct various SE tasks ranging from code translation [2], [3], commit message generation [4], [5], code review [15], [16], and program repair [6], [7] to name a few. In addition to research adoption, practitioners also integrate LLMs into software development activities. Recent years have witnessed numerous LLM-based AI code assistants, such as Github Copilot [17], JetBrains AI assistant [18], and Visual Studio IntelliCode [19]. Sergeyuk et al. [20] have found that developers tend to delegate certain SE tasks to those LLM-based AI code assistants, including test case and documentation generation.

Techniques to Improve LLM Performance. Various techniques exist to improve the performance of LLMs when automating a task. Historically, fine-tuning with domain-specific

datasets has been the go-to approach to apply the LLMs to a downstream task [21]. While fine-tuning can improve a generic LLM’s performance on the task, it requires heavy computational resources and curated datasets, making it an expensive and time-consuming approach. However, recent works have [11], [22] suggested that simply prompting the LLMs with well-designed prompt engineering techniques can be equally or more effective in certain conditions. For example, Shang et al. [23] compared fine-tuning and prompt-based methods for generating unit tests, showing that carefully structured prompts sometimes achieve comparable results. Similarly, Chen et al. [24] demonstrated that prompting LLMs outperformed fine-tuning in taxonomy-building tasks, and the performance gap widened when the LLMs were fine-tuned with limited data. Moreover, [9] offers a survey of 58 techniques, illustrating how prompt formulation can dramatically affect outcomes. This highlights the growing recognition that the way prompts are structured can significantly influence the quality of the generated output, making prompt engineering a crucial skill for developers and researchers. Despite increased interest in prompt-based methods with LLMs, best practices remain inconsistent, especially in automating software development and maintenance activities. In this study, we aim to fill this gap by providing insights into how different SE tasks would benefit from different prompting techniques so that researchers and practitioners can better utilize the power of the LLMs to automate these SE tasks.

Comparisons of LLM Approaches. Prior works have investigated the effectiveness of different pre-trained LLMs in performing SE tasks. For instance, Niu et al. [12] evaluated 19 pre-trained models on 13 SE tasks, highlighting the correlation between LLMs of different architectures and their performance on SE tasks. Similarly, a survey on code generation [25] analyzed multiple LLM-based solutions and ranked their performance across standard benchmarks. Alizadeh et al. [26] extended this line of work by analyzing the energy–accuracy trade-off of 18 LLM families, demonstrating that larger models and higher energy budgets do not always yield superior results and that no single model dominates across diverse SE tasks. While these comparisons have shown the importance of selecting the right model, they rarely revealed how the wording and structure of the prompt itself can influence the quality of the generated content, leading to an oversight, as even a well-trained model can underperform if not guided by well-constructed instructions formatted by an appropriate prompting technique [27]. It remains unclear how to adapt each task-specific scenario via prompting to maximize SE task performance, an aspect we address in this work.

In addition, we note that while both fine-tuning and prompt engineering have been studied in the domain of software engineering, no comprehensive evaluation has yet examined how different prompt engineering techniques perform across SE tasks. This comparative analysis can highlight the strengths and weaknesses of various prompting techniques, offering practical guidelines for developers and researchers alike. To fill this gap, our research systematically compares recognized

prompting techniques drawn from the literature on LLMs and SE across diverse scenarios such as bug fixing, test generation, and code summarization.

Explaining the Prompt Technique. Researchers have explored various metrics to explain model predictions and gain insights into their behavior. One such metric is *contrastive explanation*, which seeks to clarify why a model predicts a specific outcome for a given input instead of an alternative outcome. In the NLP domain, prior studies have employed *contrastive explanation* to better understand model decisions in tasks such as natural language inference [13], question answering [28], and grammatical acceptability [29]. These studies demonstrate that *contrastive explanations* can be more effective in revealing model reasoning than simply prompting the model for a rationale. In this work, we adopt the practices proposed in [13] to examine, given a prompt, task instance, and model-generated output, why one prompt technique leads to better performance than another. Our work is the first to investigate *why* certain prompts perform better in various SE tasks leveraging *contrastive explanations*.

III. METHODOLOGY

In this section, we introduce the methodology of our experiments for answering our research questions.

A. Software Engineering Tasks

To make our findings more generalizable and comprehensive, we targeted a broad range of Software Engineering (SE) tasks curated by Niu et al (Table I). [12]. These tasks have been widely explored and automated by researchers using LLMs, demonstrating various characteristics in terms of input-output modalities. They have also been grouped into code understanding and code generation tasks. Due to the limited financial budget, we could not run our experiments on all the datasets associated with these SE tasks. Thus, we conducted random sampling on each dataset with a 95% confidence level and a 5% margin of error.

The original study from which we derived our task set included a total of 13 tasks. For the purposes of our investigation, we refined this list to 10 tasks. Specifically, we excluded Code-to-Code Retrieval and Code Search as these tasks are primarily designed for code embedding models [30], [31] rather than generative models where the design of the prompting techniques can impact the performance. Additionally, we excluded the Code Completion task due to challenges in reproducing the required dataset format using the implementation provided in the replication package [12]. In total, we obtained 10 sampled datasets. Table I lists the sample size for each dataset for an SE task.

B. Prompting Techniques

A prompting technique is a blueprint that describes how to structure a prompt that serves as the input query to the LLMs. It may incorporate conditional logic, parallelism, or other architectural considerations to obtain responses from the

LLMs that better align with the user intent and task objectives [9].

Schulhoff et al. [9] curated a catalog of 46 prompting techniques in literature. To ensure that prompting techniques can be readily applied to the selected SE tasks, we devised a set of exclusion criteria to filter out those that were not relevant or applicable. This filtering process was conducted independently by two researchers, master’s students with more than two years of experience in prompt engineering. In cases where disagreement arose, a third researcher with extensive experience in LLMs and software development was consulted to resolve the divergence and reach a consensus. The whole process was conducted through a negotiated agreement [42]. These filtering criteria are: (1) Only one prompting technique was retained if multiple prompting techniques function similarly (i.e., K-Nearest Neighbor [43] and Vote-K [44] both are few-shot techniques that differ only on the method used to choose the examples); (2) the prompting technique cannot be an ensemble of multiple other techniques (i.e., DENSE [45]); (3) Since our goal was to examine the effectiveness of the prompt content and structure, the prompting technique cannot rely on external task-specific tools (i.e., ReAct [46]).

After filtering, 14 out of the 46 prompting techniques were adopted in our analysis. The list of all 46 prompting techniques and the rationale for whether they were included or not is available in our replication package ¹. We describe each of the 14 techniques as follows:

- 1) **Exemplar Selection KNN [43] (ES-KNN):** Selects exemplars using a k-nearest neighbor approach to enrich the prompt context. In detail, the data samples from the full original dataset (excluding the current data sample to be queried) are transformed into vector embeddings using a code-optimized embedding model [47]. We selected the most widely adopted and top-performing embedding model available at the time of our experiment. Depending on the SE task, since the data samples may consist of source code, natural language, or both, we selected a multilingual embedding model that was trained with both natural language and multiple programming languages. We then computed cosine similarity between the current sample and encoded data samples to retrieve the top-k nearest neighbor samples, which were selected as in-context exemplars to be added to the prompt. We selected the most frequency used and best performing
- 2) **Few Shot Contrastive CoT [48] (CCoT):** Uses both correct and incorrect chain-of-thought examples to refine reasoning steps.
- 3) **Tree Of Thought [49] (TroT):** Structures multiple branching reasoning paths to explore diverse solution strategies for complex design problems.
- 4) **Self Ask [50] (SA):** Makes the model generate its own follow-up questions before answering, which helps it break down and solve complex problems step by step.

¹<https://github.com/prompt-study/prompt-tasks-study/tree/main>

TABLE I
SUMMARY OF SELECTED SE TASKS AND EVALUATION METRICS USED

Task (Abbreviation)	Category	Dataset	Metric(s)	Sample size
Defect Detection (DD)	Code Understanding	Devign [32]	Acc	391
Clone Detection (CD)	Code Understanding	BigCloneBench [33]	F1	390
Exception Type Prediction (ET)	Code Understanding	Kanade et al. [32]	Acc	380
Code Question Answering (QA)	Code Understanding	CoSQA [34], CodeSearchNet [35]	Acc	382
Code Translation (CT)	Code Generation	CodeTrans [35]	CodeBLEU [36]	378
Bug Fixing (BF)	Code Generation	BFP [37]	CodeBLEU [36]	390
Mutant Generation (MG)	Code Generation	GM [38]	BLEU [39]	390
Assert Generation (AG)	Code Generation	ATLAS [40]	BLEU [39]	390
Code Summarization (SM)	Code Generation	DeepCom [35]	BLEU [39]	390
Code Generation (CG)	Code Generation	CONCODE [41]	CodeBLEU [36]	391

- 5) **Universal Self Consistency [51] (USC):** Combines multiple answers from the model and uses a meta-prompt to pick the most consistent one, improving output reliability.
- 6) **Self Refine [52] (SR):** Iteratively improves initial responses by self-evaluating and updating code explanations or solutions.
- 7) **Self-Generated In-Context Learning [53] (SG-ICL):** Automatically generates in-context exemplars to simulate few-shot learning, streamlining prompt formulation for coding tasks.
- 8) **Thread Of Thought [54] (ToT):** This technique guides the LLMs to work through a problem step by step, focusing on breaking down a large or complex task into smaller, manageable parts. For example, instead of solving everything at once, the model is told to pause, summarize, and analyze each step before moving to the next step, making the reasoning process more organized and clear.
- 9) **Step Back Prompting [55] (SBP):** With step back prompting, the model is first asked to examine the problem as a whole and think about the key ideas or main facts, before drafting the solution. This helps the LLMs to plan ahead and avoid jumping into details.
- 10) **Emotional Prompting [56] (EP):** Incorporates affective language to shape engaging and empathetic responses, useful in writing user-friendly documentation or error messages.
- 11) **Style Prompting [57] (SP):** Directs the model to adopt a specific tone or format, ensuring that generated code comments and documentation align with a desired style.
- 12) **Rephrase and Respond [58] (RR):** This technique asks the LLMs to first restate the question in its own words and add any extra details if needed, before giving an answer. By making the LLM explain the question to itself, the model is assumed to be more likely to fully understand what is being asked and give a more accurate and detailed response.
- 13) **Role Prompting [59] (RP):** Assigns a specific persona—such as a code reviewer or developer—to

tailor responses to the nuances of specific SE tasks.

- 14) **Analogical Prompting [60] (AP):** This technique prompts the LLMs to use analogies to make code explanations or design ideas easier to understand. This helps turn complex or abstract data structures or algorithms into relatable, real-world examples.

C. Prompt Validation

We conducted a structured prompt validation process [61] to reduce the potential impact of prompt wording on the task performance. For each prompting technique, we constructed ten prompt variation templates with different synonyms and phrasings, using OpenAI’s ChatGPT [62] as a tool to ensure no loss in text quality [63], [64] (e.g., eliminating grammatical errors, maintaining sentence coherence, and preserving the original intent of the technique), as well as to maximize variability across the generated prompts. Figure 1 contains two examples of prompt template variations. Then six researchers, who were also the authors of this study, participated in the manual process of reviewing these templates to ensure that the semantic meaning of the prompts remained consistent. To detail, a pair of researchers was assigned 4 to 5 prompting techniques. Each researcher independently reviewed the variation templates for their assigned technique to determine whether the prompts conformed to the descriptions provided by scientific literature. A variation template was accepted only if both researchers in the pair agreed. In cases of disagreement, the template was discarded, and new variation was generated and reviewed until ten acceptable template variations were finalized.

The average Cohen’s kappa for inter-rater reliability across all prompting techniques was $\kappa=0.45$.

D. Language Model Selection

To ensure the generalizability and comprehensiveness of our findings, we employed a diverse set of LLMs. These models have been developed by different organizations, focusing on various architectures and training objectives. Since the selected SE tasks involve not only natural language

"Create a variant of the code below by introducing subtle alterations. Return solely the revised code. Reflect on your modifications and trust your expertise to make a mutation that truly stands out."

Variation 1

"Please develop a modified version of the provided code by implementing minor modifications. Output only the altered code. Are you confident in your adjustments? Embrace the challenge with passion to create a subtle yet purposeful mutation."

Variation 2

Fig. 1. Example of two variations of the emotion prompting template for the task Mutant Generation.

but also source code, we chose the LLMs that show state-of-the-art performance on widely adopted code generation benchmarks, such as EvalPlus [65]. Consequently, we selected *DeepSeek-V3* [66], *Qwen2.5-Coder-32B-Instruct* [67], *Llama-3.3-70B-Instruct* [68], and *OpenAI o3-mini* [69].

E. Result Collection

The results were obtained by applying each prompting technique to all ten selected SE tasks (Table I). For each prompting technique being applied to an SE task, we divided the dataset evenly and ensured all prompt template variations were applied as close to the same number of data instances as possible (between 38 and 40 instances per variation). These prompt template variations (Section III-C) were then sent to the LLMs for solving the SE task. Specifically we used together.ai API [70] for the open source language models: *DeepSeek-V3*, *Qwen2.5-Coder-32B-Instruct* and *Llama-3.3-70B-Instruct*, and for the *o3-mini* we used the OpenAI API. The raw model responses were subsequently processed using task-specific extraction scripts. Our parser aimed to extract only the relevant answer field from the model's output according to the expected output format for each SE task. To ensure the parser did not artificially inflate performance, we only extracted answers delimited by special markers that were specified in the prompt instructions. If the model correctly used the specified delimiters, our parser would return only the content within those delimiters. If the model response failed to include the required delimiters, we included the entire response in our analysis. We included the entire model response as the answer for evaluation. This process ensured that the extraction step reflected the actual performance of the LLMs in adhering to prompt specifications and instructions. After this cleaning step, the extracted responses were used to compute each SE task's evaluation metrics (Table I).

In addition, we constructed a baseline to demonstrate the effectiveness of the prompting techniques under study. Specifically, the baseline only prompts the LLMs with a simplistic instruction to complete the SE task along with the data sample. For each of the ten SE tasks, we averaged performance metrics across the four LLMs to have an overview of the effectiveness of the prompting techniques. Then, we ranked the prompting techniques to identify the most effective one for the SE tasks.

The ranking was determined by aggregating the performance metric for each technique across four LLMs. The metrics for the SE tasks are included in Table I. We describe them in the following:

F. Linguistic Metrics

To reveal characteristics that have the potential to impact the performance of the SE tasks when prompting LLMs, we selected the following linguistic metrics that are relevant to prompt engineering:

- **Lexical Diversity (MATTR)** assesses the variety of vocabulary used in the prompt by calculating the Moving-Average Type-Token Ratio (MATTR) [71].
- **Flesch-Kincaid Grade Level** evaluates the readability of a prompt by estimating the U.S. school grade level necessary for its comprehension [72].
- **Gunning Fog Index** calculates the years of formal education a reader needs to understand the prompt [73].
- **Token Count** calculates the number of tokens in the prompt.
- **Flesch Reading Ease** rates prompt readability on a scale from 0 (very difficult) to 100 (very easy), with higher scores indicating easier-to-read prompts. This metric is widely used in software documentation and requirements engineering studies [74], [75].

We calculated these metrics for all the prompts in this study and averaged the metric values of all the prompts using a prompting technique for each of the selected SE task. The goal was to collect the linguistic characteristics of different prompting techniques for the SE tasks. To investigate their association with the prompting techniques' effectiveness in the SE tasks, we conducted Spearman correlation [76] between the averaged values of these linguistic metrics and the performance scores of the SE tasks (Table I).

G. Contrastive Explanation

To further investigate why a specific prompting technique shows better performance in an SE task than other techniques, we used contrastive explanation [13] that attempts to uncover the factors associated with the model's decisions. Contrastive explanations provide interpretable results similar to what a human would explain a decision, which inherently answers the question "*why P, rather than Q?*" [77]. Since the decisions the LLMs make are often complex, the complete explanations can be complicated and uninterpretable [78]. Contrastive explanations omit the causal attributes or factors common to both P and Q, making the explanations easy to understand [13].

To utilize contrastive explanations to explain the superiority of a prompting technique in solving the SE tasks, we prompted the LLMs using in-context learning [79]. To be more specific, for each SE task, we included the definitions of the prompting technique that performed the best and the worst (Section III-E) and prompted the LLM to generate contrastive explanations. In total, we prompted the LLM (i.e., the same LLM that had generated the responses for the SE task) to obtain comprehensive explanations regarding the potential factors that were

responsible for a prompting technique’s superior performance compared to worst prompting technique. The process was repeated for all four LLM models used in our experiment.

To ensure the quality and interpretability of the explanations, we also attached four demonstration examples (following the best practice of in-context learning in software engineering research [79]), each consisting of the queries constructed by the two prompting techniques and their responses. To make the responses of the two prompting techniques more distinguishable for the LLM to generate explanations, we only included examples whose responses by the superior prompting technique were significantly better than those of the worst prompting technique (For understanding tasks, the response of the superior prompting technique is a correct prediction while that of the compared technique is not. For generation tasks, we ensure the metric scores of the two techniques have the most difference). Due to the space limit, we include the complete prompt for all the SE tasks in our replication package. To ensure the quality of the prompt to extract explanations, we followed the best practices [80] and examined the responses after executing the composed prompts ten times.

We collected all the contrastive explanation responses provided by the four LLM models. Then, for each SE task, the responses across the models were combined for manual evaluation. To detail, we used a card sorting method [81] to conduct a qualitative analysis of these responses. Two researchers participated in this analysis. Each of them evaluated the explanations and categorized them into codes. They then met to discuss these codes, which were further grouped into high-level categories. Finally, the researchers refined the categories and organized them into meaningful themes that can sufficiently explain why a certain technique outperformed the others.

IV. RESULTS

A. (RQ1) How do different prompting techniques impact the performance of SE tasks?

Table II shows the best-performing prompting technique for each SE task using different LLMs. To have an overview of the effectiveness of the prompting techniques for each task, we also provided the prompting techniques that showed the best performance, averaged across all the LLMs in the *Aggregate* column. We note that while no single prompting technique emerged as universally optimal across all tasks, there are certain prompting techniques that tend to show superior effectiveness than others. For example, for tasks such as *Clone Detection*, *Code Translation*, and *Assert Generation*, *ES-KNN* consistently outperforms other prompting techniques across all four selected LLMs. This suggests that semantically similar in-context demonstration examples in the prompt may be able to hint the LLMs how to process the current data sample to achieve the goals of the SE tasks properly. In addition, *USC* demonstrates the best performance in *Code QA* and *Code Generation* for most of the LLMs, which indicates that providing the model with structured examples or encouraging exploration of multiple solution pathways can mitigate errors

in the generated output. It is worth noting that *ToT* presents outstanding performance for the task of *Defect Detection*. It prompts the LLMs to engage in a process where the code is broken down into components, and each component is examined step by step to identify the potential defects. Such a process may push the LLMs to focus on specific components one at a time, increasing the chance of identifying defects.

As for the four selected LLMs, *DeepSeek-V3*, *Qwen2.5-Coder-32B-Instruct*, and *Llama-3.3-70B* show nearly consistent results regarding the prompting technique that is the most effective across all SE tasks (i.e., *ES-KNN* is the most effective one across all these three models in *Clone Detection*, *Exception Type Prediction*, *Code Translation*, and *Assert Generation*). However, the outperforming prompting techniques that are applied to *OpenAI o3-mini* are different for most SE tasks, with *RP* being the majority winner.

We also include the worst-performing prompting techniques for the SE tasks in Table III. Table IV lists the prompting techniques along with the average values of the evaluation metrics across the four LLMs for the ten SE tasks in this study. We also included the aggregate performance of the two task groups: code understanding and generation tasks, which are displayed in z-score [82] format to make comparing different metrics possible. Here, a z-score indicates how many standard deviations a particular value is from the mean, allowing for a standardized comparison of performance across different metrics and scales.

It shows that the worst-performing prompting techniques even underperform the baseline (Section III-E) where simplistic instruction is used. This finding points out that complex prompting techniques would not necessarily improve performance when applied to achieve certain goals. Our results reveal compatibilities of various prompting techniques to different SE tasks, providing practical guidelines to researchers and software developers. These guidelines can be essential to ensure software quality since LLMs are increasingly being used to automate software development and maintenance activities [20].

Observation 1: No single prompt technique consistently outperforms others across all SE tasks. Certain prompting techniques can negatively impact their performance.

B. (RQ2) What linguistic features of prompting techniques are associated with improved performance on SE tasks?

To examine how linguistic characteristics of prompts relate to task performance, we analyzed selected linguistic metrics (Section III-F) across various task groups. Note that the tasks are grouped into code understanding tasks and code generation tasks following [12] (Table I). Our goal was to identify potential associations between these features and SE task performance. The results are summarized below:

Lexical diversity, measured by the Moving-Average Type-Token Ratio (MATTR), shows a strong positive correlation with performance across all tasks (aggregate: $r = 0.4440$, $p < 0.001$). This correlation is particularly pronounced in code

TABLE II
BEST PROMPTING TECHNIQUES FOR EACH MODEL AND AGGREGATE

Task	Aggregate	Qwen	DeepSeek	Llama	o3-mini
Understanding Tasks	ES-KNN	ES-KNN	ES-KNN	ES-KNN	RP
Defect detection	ToT	ToT	ToT	ToT	RP
Clone detection	ES-KNN	ES-KNN	ES-KNN	ES-KNN	ES-KNN
Exception type	ES-KNN	ES-KNN	ES-KNN	ES-KNN	RR
Code QA	USC	SG-ICL	USC	USC	RP
Generation Tasks	ES-KNN	ES-KNN	ES-KNN	ES-KNN	ES-KNN
Code translation	ES-KNN	ES-KNN	ES-KNN	ES-KNN	ES-KNN
Bug fixing	Control	SG-ICL	SG-ICL	SG-ICL	Control
Mutant generation	ES-KNN	RP	ES-KNN	ES-KNN	RP
Assert generation	ES-KNN	ES-KNN	ES-KNN	ES-KNN	ES-KNN
Code summarization	Control	SG-ICL	ES-KNN	Control	SG-ICL
Code generation	USC	USC	SG-ICL	USC	USC

TABLE III
WORST PROMPTING TECHNIQUES FOR EACH MODEL AND AGGREGATE

Task	Aggregate	Qwen	DeepSeek	Llama	o3-mini
Understanding Tasks	SR	SR	CCoT	SR	SBP
Defect detection	EP	RR	EP	EP	SBP
Clone detection	SR	SG-ICL	SG-ICL	USC	SR
Exception type	RR	SR	ToT	SR	SR
Code QA	TroT	TroT	TroT	Control	Control
Generation Tasks	SG-ICL	SG-ICL	ToT	SG-ICL	SR
Code translation	SG-ICL	SG-ICL	USC	TroT	USC
Bug fixing	RR	ToT	RR	SR	SA
Mutant generation	ToT	ToT	SG-ICL	USC	SR
Assert generation	SR	TroT	SG-ICL	USC	SR
Code summarization	RR	RR	RR	SA	RR
Code generation	ES-KNN	TroT	TroT	TroT	ES-KNN

generation tasks ($r = 0.5229$, $p < 0.001$) and also significant in code understanding tasks ($r = 0.3468$, $p = 0.0088$). These results suggest that prompts with a rich vocabulary and a diverse use of words can enhance model performance.

Token count exhibits a negative correlation with performance in both code understanding ($r = -0.2567$, $p = 0.0022$) and code generation tasks ($r = -0.3200$, $p = 0.0030$). This indicates that longer prompts may not necessarily lead to better outcomes.

Readability metrics show varying correlations depending on the SE tasks. In code understanding tasks, Flesch-Kincaid Grade Level scores correlate negatively with task performance ($r = -0.2974$, $p = 0.0260$), suggesting that simpler prompts are more effective. However, in code generation tasks, there’s a positive correlation ($r = 0.2975$, $p = 0.0060$), implying that more complex prompts may benefit these tasks.

Similarly, the Gunning Fog Index shows a negative correlation in code understanding tasks ($r = -0.3795$, $p = 0.0039$) and a positive correlation in code generation tasks ($r = 0.2938$, $p = 0.0067$).

However, The Flesch Reading Ease score shows different correlations. It correlates positively with task performance in code understanding tasks ($r = 0.2797$, $p = 0.0369$) and negatively with code generation tasks ($r = -0.3366$,

$p = 0.0017$). This would indicate that readability impacts performance differently depending on the task characteristics, but different metrics demonstrate different correlation patterns. We argue that future research should further investigate how the readability of the prompts impacts the performance of the LLMs in the SE tasks. Therefore, our results underscore the importance of tailoring prompt characteristics to the specific nature of the task. Enhancing lexical diversity appears universally beneficial, while optimizing readability requires further investigation.

Observation 2: Linguistic features of prompts exhibit significant and various correlations with SE task performance.

C. (RQ3) What factors, according to LLMs, contribute to the effectiveness of prompting techniques in SE tasks

To understand why a specific prompting technique outperforms others in SE tasks, we used contrastive explanations to examine the additional factors contributing to the effectiveness of each technique. Table IV presents the categorized factors of the best prompting technique in contrast to the worst techniques, which was manually evaluated as outlined in III-G. For example, in the code translation task, the best prompting technique (*ES-KNN*) included structured guidance and in-context examples, in contrast to the worst technique (Three-of-Thought). Categories of factors in Table IV represent additional factors incorporated into the best prompting technique that contributed to superior task performance. The findings from the contrastive explanation analysis suggest that *Structured Guidance* and *In-Context Examples* are the most prevalent factors among the best-performing techniques. Furthermore, these common factors consistently appear across the majority of the SE tasks included in our experiment. This indicates that, for most SE tasks, effective prompts that enhance task performance tend to provide the model with structured guidance and relevant in-context examples. Below is the list of all categories for the factors identified for the best prompting techniques based on the manual analysis explained in Section III-G

- **Ambiguity Reduction (15.38%):** Instructions to clarify task requirements.
- **Correctness and Precision (2.56%):** Instructions to ensure relevant and accurate task output.
- **Efficiency (17.95%):** Instructions to deduce an optimal solution.
- **In-Context Example (21.80%):** Contains examples relevant to the task.
- **Reasoning (2.56%):** Instructions to include detailed justifications to support the response.
- **Robustness/Comprehensiveness (7.70%):** Instructions to ensure implementations satisfy all task-specific requirements.
- **Structured Guidance (32.05%):** Includes instructions to follow the conventions and patterns aligned with task output.

TABLE IV
COMPARISON OF PROMPTING TECHNIQUES WITH ASSOCIATED REASON CATEGORIES

Task	Std	Best		Control Prompt		Worst		Best vs Worst Contrastive Explanation
		Technique	z-score	z-score		Technique	z-score	
Understanding Tasks		ES-KNN	1.12	0.29		ToT	-0.66	
Defect detection	3.60	ToT	73.66	66.22		SR	59.91	Structured Guidance, Robustness
Clone detection	1.17	ES-KNN	68.60	66.03		SG-ICL	63.45	Structured Guidance, Efficiency
Exception type	7.16	ES-KNN	82.50	78.16		ToT	52.04	Structured Guidance, In-Context Example
Code QA	1.39	USC	55.67	50.99		TroT	51.44	Ambiguity Reduction, Structured Guidance
Generation Tasks		ES-KNN	1.83	0.80		TroT	-2.11	
Code translation	6.30	ES-KNN	42.08	30.19		TroT	13.39	Structured Guidance, In-Context Example
Bug fixing	4.69	SG-ICL	36.02	36.26		TroT	16.45	In-Context Example, Robustness
Mutant generation	19.59	ES-KNN	69.93	67.98		TroT	16.44	Structured Guidance, In-Context Example
Assert generation	15.23	ES-KNN	65.44	25.24		TroT	0.92	Structured Guidance, In-Context Example
Code summarization	1.51	SG-ICL	4.15	4.16		ToT	0.45	Structured Guidance, In-Context Example
Code generation	3.04	USC	24.44	23.18		TroT	13.45	Ambiguity Reduction, Efficiency

Observation 3: Prompting techniques show greater effectiveness when they include structured guidance aligned with task objectives and relevant examples.

D. (RQ4) How are resource costs associated with the performance of prompting techniques in SE tasks?

To investigate how the resource costs impact the performance of the prompting techniques, we examine how resource-efficient each prompting technique is in conducting the SE tasks. To do so, we considered two resource measurements, namely, the number of tokens in the prompt and the response time. We normalize the efficiency of each prompting technique by dividing the values of the performance evaluation metrics by the number of used tokens and the response time. Finally, the prompting techniques are ranked based on the resource-efficiency. We present the most and least resource-efficient prompting techniques in Table VII and VIII, along with the original best and worst prompting techniques in Section IV-A.

Our findings reveal that, while *ES-KNN* remains at the top for many tasks in Section IV-A, it is not the most token-efficient one. The additional context and examples provided by *ES-KNN* increase token consumption, making it less suitable in settings where minimizing token usage is critical. In contrast, *RP* consistently ranks at or near the top for token efficiency across most tasks.

Table V presents the mean number of tokens saved per prompt and an aggregate average. Code generation shows the highest token savings overall (Aggregate: 10,733.42), especially for Llama (8837.50), Qwen (9278.82), and DeepSeek (8502.98). Bug fixing also has very high token savings (Aggregate: 8306.71), with DeepSeek (8744.93), Llama (7621.01), and Qwen (7698.30). However, o3-mini lags far behind (360.08). Code QA achieves an exceptionally high savings with Llama (19485.09), much higher than any other model on any task, which may indicate unusually long original prompts or strong compressibility. o3-mini consistently shows the low-

TABLE V
MEAN NUMBER OF TOKENS SAVED PER PROMPT

Task	Aggregate	Qwen	DeepSeek	Llama	o3-mini
Defect detection	2361.47	2195.55	3079.72	4867.86	2559.22
Clone detection	3588.50	3518.13	3777.58	3498.83	3634.21
Exception type	1985.42	1957.17	2070.51	1913.72	2583.11
Code QA	2143.98	2812.42	3984.25	19485.09	1696.93
Code translation	3021.62	3094.46	3265.77	3471.79	551.83
Bug fixing	8306.71	7698.30	8744.93	7621.01	360.08
Mutant generation	2945.01	92.55	2679.70	4458.38	362.61
Assert generation	250.75	427.77	3886.71	5269.65	643.66
Code summarization	3129.23	3408.88	4586.85	2648.01	1004.58
Code generation	10733.42	9278.82	8502.98	8837.50	791.71

TABLE VI
MEAN VALUE OF TIME (SECONDS) SAVED PER PROMPT

Task	Aggregate	Qwen	DeepSeek	Llama	o3-mini
Defect detection	107.60	220.09	103.48	56.87	66.87
Clone detection	80.31	191.34	83.02	30.81	32.45
Exception type	66.68	186.78	22.89	19.91	54.61
Code QA	85.73	207.76	68.15	39.70	43.34
Code translation	111.33	101.74	107.54	39.80	4.41
Bug fixing	211.72	363.59	262.24	72.23	2.80
Mutant generation	114.80	17.89	91.95	48.83	2.54
Assert generation	139.91	91.82	129.99	55.94	5.06
Code summarization	121.34	269.99	143.31	32.81	56.34
Code generation	231.05	384.03	240.96	78.14	7.55

est token savings across most tasks, suggesting its prompts may already be concise or less responsive to compression. Tasks with lower overall savings include Assert generation (Aggregate: 250.75), where DeepSeek and Llama still achieve large reductions (3886.71 and 5269.65, respectively), but Qwen and o3-mini do not. Mutant generation also shows lower savings (Aggregate: 2945.01), where Qwen performs particularly poorly (92.55).

A similar shift occurs when response time is considered. Techniques such as *ES-KNN* and *SA* are frequently among the fastest, while approaches like *USC* and *SR* consistently fall to the bottom due to the computational overhead associated with their iterative reasoning or the generation of multiple responses. Despite performance gain in several code generation tasks such as *Code QA* and *Code Generation*, they tend to incur additional time costs to obtain the responses.

Table VI presents the average number of seconds saved per prompt. Similar to token savings, Code Generation (231.05s) and Bug Fixing (211.72s) achieved the largest average time savings across models, reinforcing that longer and more complex tasks benefit most from prompt compression. The tasks that benefit least from prompting across models are Exception type prediction (66.68s), Clone detection (80.31s), and Code QA (85.73s). This suggests that prompting techniques provide less efficiency gain for tasks that may be more recognition- or classification-based, rather than generative in nature.

Observation 4: Prompting techniques like ES-KNN deliver the highest performance and speed but include more tokens in the prompts, while techniques such as USC boost performance but consume more time to obtain a response.

V. DISCUSSION

Our results show that prompting effectiveness varies by task and LLM, with no single technique consistently outperforming others. Exemplar Selection KNN was strong across tasks, while Role Prompting offered a cost-efficient alternative. Notably, the *o3-mini* model showed different performance patterns. These findings have several implications. First, they highlight the importance of prompt-model alignment, suggesting that prompt engineering should be adaptive, not one-size-fits-all. Practitioners should empirically validate prompting strategies when switching models, even within the same model family. Second, the variation in performance across tasks reinforces the need for task-aware prompting, where the design of the prompt accounts for the structure, complexity, and requirements of the underlying task. Finally, the deviations seen with *o3-mini* raise concerns about prompt robustness across model scales. This suggests that lighter-weight models may require distinct prompting techniques or additional tuning to match the effectiveness seen in larger models. These insights emphasize the value of prompt evaluation frameworks that account for model-task combinations and point to future research opportunities in automated prompt selection methods.

Our contrastive explanation analysis revealed that prompting techniques like Exemplar Selection KNN outperform others by offering clear structural guidance and relevant in-context examples, enhancing performance across most SE tasks. These findings indicate that effective prompt design requires more than wording—it demands well-structured guidance and tailored examples aligned with the task. Prioritizing prompt clarity and minimizing ambiguity can boost model robustness

and efficiency, helping practitioners develop more effective prompting strategies to maximize LLM performance in diverse software engineering tasks. Moreover, these insights suggest that adaptive prompt engineering, which dynamically incorporates task-specific structures and examples, could further improve results. This also highlights the potential for automated prompt optimization tools to assist developers in crafting prompts that align with task requirements, ultimately accelerating development workflows and reducing trial-and-error in prompt design.

Our analysis of token savings and time savings reveals that prompt optimization offers significant efficiency gains across SE tasks, with clear variation based on task type and LLM architecture. Tasks such as Code Generation and Bug Fixing exhibited the highest average token savings, exceeding 8,000 tokens per prompt. This suggests that these token-intensive tasks benefit substantially from prompt compression strategies, with potential reductions in inference cost and latency. These are critical considerations in real-world applications that utilize commercial LLMs with token-based billing.

We also observed considerable variation in token and time savings across LLMs. Larger models like Llama and DeepSeek consistently achieved higher savings, indicating they are more amenable to prompt compression, possibly due to differences in tokenization schemes, context window capacity, or internal architectural design. In contrast, *o3-mini* showed relatively limited token savings across most tasks, suggesting a constrained capacity for compression. These findings reinforce that prompt optimization strategies should be model-aware, as techniques effective on one LLM may not generalize to another.

Task-specific trends also emerged. For example, Assert Generation and Exception Type classification yielded lower aggregate token savings, implying that prompts for these tasks may already be concise, or that further compression risks omitting critical context. Conversely, Code QA with Llama showed exceptionally high savings, warranting further investigation into whether such tasks include redundancies or patterns that facilitate aggressive but safe compression. Some tasks, like Exception Type and Clone Detection, consistently had lower time savings across models, suggesting their input prompts may already be compact or that compression has limited impact on runtime.

From a practical standpoint, these findings emphasize that prompt design should carefully balance efficiency and performance. While reducing token usage is desirable, it must not come at the expense of output quality, particularly for complex tasks like Bug Fixing and Code Translation, which may require more detailed and nuanced input. The observed differences across tasks and models highlight an opportunity for automated, context-aware prompt rewriting tools that tailor prompt structure dynamically to maximize utility and efficiency.

VI. THREATS TO VALIDITY

In this section, we outline the potential threats to the validity of our study.

TABLE VII
RESOURCE COSTS ASSOCIATED WITH BEST PERFORMING TECHNIQUES (DEFAULT / TOKEN / TIME)

Task	Aggregate	Qwen	DeepSeek	Llama	o3-mini
Understanding Tasks	ES-KNN / Control / ES-KNN	ES-KNN / Control / Control	ES-KNN / Control / SA	ES-KNN / RP / ES-KNN	RP / Control / Control
Defect detection	ToT / Control / ES-KNN	ToT / Control / Control	ToT / Control / SA	ToT / RP / ES-KNN	RP / Control / Control
Clone detection	ES-KNN / SP / ES-KNN	ES-KNN / RP / Control	ES-KNN / EP / TroT	ES-KNN / EP / ES-KNN	ES-KNN / Control / Control
Exception type	ES-KNN / Control / ES-KNN	ES-KNN / Control / ES-KNN	ES-KNN / RP / TroT	ES-KNN / Control / ES-KNN	RR / Control / Control
Code QA	USC / Control / ES-KNN	SG-ICL / Control / ES-KNN	USC / Control / ES-KNN	USC / RP / ES-KNN	RP / Control / Control
Generation Tasks	ES-KNN / Control / ES-KNN	ES-KNN / Control / ES-KNN	ES-KNN / RP / SA	ES-KNN / Control / ES-KNN	ES-KNN / Control / Control
Code translation	ES-KNN / Control / ES-KNN	ES-KNN / Control / ES-KNN	ES-KNN / SP / SA	ES-KNN / Control / ES-KNN	ES-KNN / Control / Control
Bug fixing	Control / Control / ES-KNN	SG-ICL / Control / Control	SG-ICL / RP / TroT	SG-ICL / Control / ES-KNN	Control / Control / Control
Mutant generation	ES-KNN / Control / ES-KNN	RP / Control / ES-KNN	ES-KNN / RP / SA	ES-KNN / Control / ES-KNN	RP / Control / Control
Assert generation	ES-KNN / Control / ES-KNN	ES-KNN / Control / ES-KNN	ES-KNN / Control / ES-KNN	ES-KNN / Control / ES-KNN	ES-KNN / Control / Control
Code summarization	Control / Control / Control	SG-ICL / Control / Control	ES-KNN / Control / ES-KNN	Control / Control / Control	SG-ICL / Control / Control
Code generation	USC / Control / ES-KNN	USC / Control / ES-KNN	SG-ICL / SA / ES-KNN	USC / Control / ES-KNN	USC / Control / Control

TABLE VIII
RESOURCE COSTS ASSOCIATED WITH WORST PERFORMING TECHNIQUES (DEFAULT / TOKEN / TIME)

Task	Aggregate	Qwen	DeepSeek	Llama	o3-mini
Understanding Tasks	SR / ES-KNN / SR	SR / SR / SR	CCoT / USC / ToT	SR / USC / SR	SBP / SG-ICL / SR
Defect detection	EP / SR / SR	RR / SBP / SR	EP / SR / SBP	EP / ES-KNN / SR	SBP / SBP / SR
Clone detection	SR / SR / SBP	SG-ICL / SR / SR	SG-ICL / USC / SBP	USC / USC / SR	SR / SR / SBP
Exception type	RR / SBP / SBP	SR / SBP / SR	ToT / SBP / SBP	SR / SR / SR	SR / SBP / SR
Code QA	TroT / SR / SR	TroT / USC / SR	TroT / USC / SR	Control / TroT / USC	Control / SG-ICL / SG-ICL
Generation Tasks	SG-ICL / SBP / SG-ICL	SG-ICL / TroT / SR	ToT / SBP / SBP	SG-ICL / SG-ICL / SG-ICL	SR / USC / USC
Code translation	SG-ICL / SBP / SG-ICL	SG-ICL / SG-ICL / USC	USC / SBP / SBP	TroT / SG-ICL / SG-ICL	USC / USC / USC
Bug fixing	RR / SR / SR	ToT / SBP / SR	RR / SR / SR	SR / SR / SR	SA / USC / USC
Mutant generation	ToT / SG-ICL / SG-ICL	ToT / USC / USC	SG-ICL / SBP / ToT	USC / SG-ICL / SG-ICL	SR / SR / SR
Assert generation	SR / SG-ICL / SG-ICL	TroT / TroT / TroT	SG-ICL / SG-ICL / SG-ICL	USC / TroT / SR	SR / SR / SR
Code summarization	RR / ToT / ToT	RR / ToT / ToT	RR / ToT / ToT	SA / ToT / ToT	RR / USC / ToT
Code generation	ES-KNN / SR / SR	TroT / SR / SR	TroT / SR / SR	TroT / TroT / SR	ES-KNN / USC / USC

Construct Validity: To mitigate bias in manual evaluation for prompting technique selection (Section III-B and prompt variation template selection (Section III-C), the researchers independently conducted the manual analysis. Following open-coding practices [83], discrepancies were resolved through negotiated agreement. Moreover, to ensure the quality of model-generated responses through contrastive explanations, we adhered to established prompting practices [13].

Internal Validity: To ensure that task performance reflects solely the effectiveness of the prompting technique rather than prompt wording or model temperature, we used a fixed default temperature value (1 for o3-mini and Deepseek, 0.7 for Llama and Qwen) for all LLM and we employed the prompt validation process as outlined in Section III-C. In addition, we adopted the dataset used by a prior study [9]. We used the same sampled dataset for each SE task to compare the performance of different prompting techniques across four LLM models.

External Validity: This study uses four LLMs to investigate 10 software engineering (SE) tasks and 14 prompting techniques. While our findings lay a foundation for the systematic evaluation of prompting techniques across diverse SE tasks, they may not be fully generalizable beyond the specific tasks, datasets, prompting techniques, and LLMs used in this work. To minimize the threat pertaining to LLMs, we selected models with distinct training objectives (general-

purpose vs. code-specialized) and model availability (open-source vs. proprietary models) and strong performance on the widely-used EvalPlus leaderboard [65].

VII. CONCLUSION AND FUTURE WORK

In this study, we conducted systematic evaluation of 14 prompt engineering techniques across 10 SE tasks using four LLMs. Our results offer concrete insights into *which* techniques yield highest performance gains and *where* resource overheads may present practical limitations in real-world deployments. The findings of our work provide empirical evidence to guide practitioners and researchers for selecting optimal prompting techniques that are best aligned with task-specific objectives and real-world operational constraints such as execution time and token usage.

Future work should extend beyond the linguistic characteristics of natural language prompts to explore additional dimensions, such as properties of prompt embeddings. Further research is needed to investigate strategies to optimize the prompt structure in accordance with task-specific requirements and dataset characteristics. All research artifacts are available on our companion website. [84]

REFERENCES

- [1] A. Vaswani and et al., "Attention is all you need," in *Advances in Neural Information Processing Systems*, vol. 30, 2017.

- [2] R. Pan, A. R. Ibrahimzada, R. Krishna, D. Sankar, L. P. Wassi, M. Merler, B. Sobolev, R. Pavuluri, S. Sinha, and R. Jabbarvand, "Lost in translation: A study of bugs introduced by large language models while translating code," in *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*, 2024, pp. 1–13.
- [3] M. Wu, J. Xu, and L. Wang, "Transagents: Build your translation company with language agents," in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, 2024, pp. 131–141.
- [4] J. Li, D. Faragó, C. Petrov, and I. Ahmed, "Only diff is not enough: Generating commit messages leveraging reasoning and action of large language model," *Proceedings of the ACM on Software Engineering*, vol. 1, no. FSE, pp. 745–766, 2024.
- [5] —, "Consider what humans consider: Optimizing commit message leveraging contexts considered by human," *arXiv preprint arXiv:2503.11960*, 2025.
- [6] S. Hou, Y. Liu, and J. Lee, "Towards an understanding of large language models in software engineering tasks," in *42nd International Conference on Software Engineering (ICSE)*, 2024, pp. 754–768.
- [7] J. Xu, Y. Fu, S. H. Tan, and P. He, "Aligning the objective of llm-based program repair," *arXiv preprint arXiv:2404.08877*, 2024.
- [8] M. Sclar, Y. Choi, Y. Tsvetkov, and A. Suhr, "Quantifying language models' sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting," in *International Conference on Learning Representations (ICLR)*, 2024.
- [9] S. S. et al., "The prompt report: A systematic survey of prompt engineering techniques," 2025. [Online]. Available: <https://arxiv.org/abs/2406.06608>
- [10] P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, and G. Neubig, "Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing," *ACM Comput. Surv.*, vol. 55, no. 9, Jan. 2023. [Online]. Available: <https://doi.org/10.1145/3560815>
- [11] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. H. Chi, Q. V. Le, and D. Zhou, "Chain-of-thought prompting elicits reasoning in large language models," in *Proceedings of the 36th International Conference on Neural Information Processing Systems*, ser. NIPS '22. Red Hook, NY, USA: Curran Associates Inc., 2022.
- [12] C. Niu, C. Li, V. Ng, D. Chen, J. Ge, and B. Luo, "An empirical comparison of pre-trained models of source code," in *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 2023, pp. 2136–2148.
- [13] A. Jacovi, S. Swayamdipta, S. Ravfogel, Y. Elazar, Y. Choi, and Y. Goldberg, "Contrastive explanations for model interpretability," *arXiv preprint arXiv:2103.01378*, 2021.
- [14] K. Alizadeh, S. I. Mirzadeh, D. Belenko, S. Khatamifard, M. Cho, C. C. Del Mundo, M. Rastegari, and M. Farajtabar, "Llm in a flash: Efficient large language model inference with limited memory," pp. 12 562–12 584, 2024.
- [15] Q. Guo, J. Cao, X. Xie, S. Liu, X. Li, B. Chen, and X. Peng, "Exploring the potential of chatgpt in automated code refinement: An empirical study," in *Proceedings of the 46th IEEE/ACM International Conference on Software Engineering*, 2024, pp. 1–13.
- [16] O. B. Sghaier, M. Weyssow, and H. Sahraoui, "Harnessing large language models for curated code reviews," *arXiv preprint arXiv:2502.03425*, 2025.
- [17] G. Copilot, "Ai that builds with you," 2023. [Online]. Available: <https://github.com/features/copilot>
- [18] J. AI, "Jetbrains ai: Optimize your workflow. with ai built for you," 2024. [Online]. Available: <https://www.jetbrains.com/ai/>
- [19] V. S. IntelliCode, "Jtype less, code more. visual studio intellcode brings ai assistance directly into your personal development flow," 2024. [Online]. Available: <https://visualstudio.microsoft.com/services/intellcode/>
- [20] A. Sergeyuk, Y. Golubev, T. Bryksin, and I. Ahmed, "Using ai-based coding assistants in practice: State of affairs, perceptions, and ways forward," *Information and Software Technology*, vol. 178, p. 107610, 2025.
- [21] B. Hanindhito, B. Patel, and L. K. John, "Large language model fine-tuning with low-rank adaptation: A performance exploration," in *Proceedings of the 16th ACM/SPEC International Conference on Performance Engineering (ICPE '25)*. Toronto, ON, Canada: ACM, 2025, pp. 92–104. [Online]. Available: <https://doi.org/10.1145/3676151.3719377>
- [22] T. B. e. a. Brown, "Language models are few-shot learners," in *Proceedings of the 34th International Conference on Neural Information Processing Systems*, ser. NIPS '20. Curran Associates Inc., 2020.
- [23] Y. Shang, Q. Zhang, C. Fang, S. Gu, J. Zhou, and Z. Chen, "A large-scale empirical study on fine-tuning large language models for unit testing," 2024. [Online]. Available: <https://arxiv.org/abs/2412.16620>
- [24] B. Chen, F. Yi, and D. Varró, "Prompting or fine-tuning? a comparative study of large language models for taxonomy construction," in *2023 ACM/IEEE International Conference on Model Driven Engineering Languages and Systems Companion (MODELS-C)*. IEEE, 2023, pp. 588–596.
- [25] J. Jiang, F. Wang, J. Shen, S. Kim, and S. Kim, "A survey on large language models for code generation," *arXiv preprint arXiv:2406.00515*, 2024.
- [26] N. Alizadeh, B. Belchev, N. Saurabh, P. Kelbert, and F. Castor, "Language models in software development tasks: An experimental analysis of energy and accuracy," 2025. [Online]. Available: <https://arxiv.org/abs/2412.00329>
- [27] P. Hase, I. Purohit, and M. Bansal, "Prompting is programming: A query language for large language models," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2023. [Online]. Available: <https://openreview.net/forum?id=RIu5lyNXjT>
- [28] Q. Chen, G. Xu, M. Yan, J. Zhang, F. Huang, L. Si, and Y. Zhang, "Distinguish before answer: Generating contrastive explanation as knowledge for commonsense question answering," *arXiv preprint arXiv:2305.08135*, 2023.
- [29] K. Yin and G. Neubig, "Interpreting language models with contrastive explanations," *arXiv preprint arXiv:2202.10419*, 2022.
- [30] Y. Liu, R. Meng, S. Joty, S. Savarese, C. Xiong, Y. Zhou, and S. Yavuz, "Codexembed: A generalist embedding model family for multilingual and multi-task code retrieval," *arXiv preprint arXiv:2411.12644*, 2024.
- [31] V. AI, "voyage-code-2: Elevate your code retrieval," 2024. [Online]. Available: <https://blog.voyageai.com/2024/01/23/voyage-code-2-elevate-your-code-retrieval/>
- [32] A. Kanade, P. Maniatis, G. Balakrishnan, and K. Shi, "Learning and evaluating contextual embedding of source code," in *Proceedings of the 37th International Conference on Machine Learning*, ser. ICML'20. JMLR.org, 2020.
- [33] J. Svajlenko, J. F. Islam, I. Keivanloo, C. K. Roy, and M. M. Mia, "Towards a big data curated benchmark of inter-project code clones," in *2014 IEEE International Conference on Software Maintenance and Evolution*, 2014, pp. 476–480.
- [34] J. Huang, D. Tang, L. Shou, M. Gong, K. Xu, D. Jiang, M. Zhou, and N. Duan, "CoSQA: 20,000+ web queries for code search and question answering," in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, C. Zong, F. Xia, W. Li, and R. Navigli, Eds. Online: Association for Computational Linguistics, Aug. 2021, pp. 5690–5700. [Online]. Available: <https://aclanthology.org/2021.acl-long.442/>
- [35] S. L. et al., "Codexglue: A machine learning benchmark dataset for code understanding and generation," 2021. [Online]. Available: <https://arxiv.org/abs/2102.04664>
- [36] S. Ren, D. Guo, S. Lu, L. Zhou, S. Liu, D. Tang, N. Sundaresan, M. Zhou, A. Blanco, and S. Ma, "Codebleu: a method for automatic evaluation of code synthesis," 2020. [Online]. Available: <https://arxiv.org/abs/2009.10297>
- [37] M. Tufano, C. Watson, G. Bavota, M. D. Penta, M. White, and D. Poshvanyk, "An empirical study on learning bug-fixing patches in the wild via neural machine translation," *ACM Transactions on Software Engineering and Methodology (TOSEM)*, vol. 28, no. 4, pp. 1–29, 2019.
- [38] G. Ye, T. Hu, Z. Tang, Z. Fan, S. H. Tan, B. Zhang, W. Qian, and Z. Wang, "A generative and mutational approach for synthesizing bug-exposing test cases to guide compiler fuzzing," in *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2023, pp. 1127–1139.
- [39] K. Papineni, S. Roukos, T. Ward, and W. jing Zhu, "Bleu: a method for automatic evaluation of machine translation," 2002, pp. 311–318.
- [40] C. Watson, M. Tufano, K. Moran, G. Bavota, and D. Poshvanyk, "On learning meaningful assert statements for unit test cases," in *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*, ser. ICSE '20. New York, NY, USA: Association for Computing Machinery, 2020, p. 1398–1409. [Online]. Available: <https://doi.org/10.1145/3377811.3380429>

- [41] S. Iyer, I. Konstas, A. Cheung, and L. Zettlemoyer, "Mapping language to code in programmatic context," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, E. Riloff, D. Chiang, J. Hockenmaier, and J. Tsujii, Eds. Brussels, Belgium: Association for Computational Linguistics, Oct.-Nov. 2018, pp. 1643–1652. [Online]. Available: <https://aclanthology.org/D18-1192/>
- [42] J. Forman and L. Damschroder, "Qualitative content analysis," in *Empirical methods for bioethics: A primer*. Emerald Group Publishing Limited, 2007, pp. 39–62.
- [43] J. Liu, D. Shen, Y. Zhang, B. Dolan, L. Carin, and W. Chen, "What makes good in-context examples for GPT-3?" in *Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures*, E. Agirre, M. Apidianaki, and I. Vulić, Eds. Dublin, Ireland and Online: Association for Computational Linguistics, May 2022, pp. 100–114. [Online]. Available: <https://aclanthology.org/2022.deeLIO-1.10/>
- [44] H. Su, J. Kasai, C. H. Wu, W. Shi, T. Wang, J. Xin, R. Zhang, M. Ostendorf, L. Zettlemoyer, N. A. Smith, and T. Yu, "Selective annotation makes language models better few-shot learners," 2022. [Online]. Available: <https://arxiv.org/abs/2209.01975>
- [45] M. Khalifa, L. Logeswaran, M. Lee, H. Lee, and L. Wang, "Exploring demonstration ensembling for in-context learning," *arXiv preprint arXiv:2308.08780*, 2023.
- [46] S. Yao, J. Zhao, D. Yu, N. Du, I. Shafraan, K. Narasimhan, and Y. Cao, "React: Synergizing reasoning and acting in language models," in *International Conference on Learning Representations (ICLR)*, 2023.
- [47] H. Face, "jina-embeddings-v2-base-code," 2025. [Online]. Available: <https://huggingface.co/jinaai/jina-embeddings-v2-base-code>
- [48] Y. K. Chia, G. Chen, L. A. Tuan, S. Poria, and L. Bing, "Contrastive chain-of-thought prompting," 2023. [Online]. Available: <https://arxiv.org/abs/2311.09277>
- [49] S. Yao, D. Yu, J. Zhao, I. Shafraan, T. L. Griffiths, Y. Cao, and K. Narasimhan, "Tree of thoughts: Deliberate problem solving with large language models," 2023. [Online]. Available: <https://arxiv.org/abs/2305.10601>
- [50] O. Press, M. Zhang, S. Min, L. Schmidt, N. A. Smith, and M. Lewis, "Measuring and narrowing the compositionality gap in language models," 2023. [Online]. Available: <https://arxiv.org/abs/2210.03350>
- [51] X. Chen, R. Aksitov, U. Alon, J. Ren, K. Xiao, P. Yin, S. Prakash, C. Sutton, X. Wang, and D. Zhou, "Universal self-consistency for large language model generation," 2023. [Online]. Available: <https://arxiv.org/abs/2311.17311>
- [52] A. M. et al., "Self-refine: Iterative refinement with self-feedback," 2023. [Online]. Available: <https://arxiv.org/abs/2303.17651>
- [53] H. J. Kim, H. Cho, J. Kim, T. Kim, K. M. Yoo, and S. goo Lee, "Self-generated in-context learning: Leveraging auto-regressive language models as a demonstration generator," 2022. [Online]. Available: <https://arxiv.org/abs/2206.08082>
- [54] Y. Zhou, X. Geng, T. Shen, C. Tao, G. Long, J.-G. Lou, and J. Shen, "Thread of thought unraveling chaotic contexts," 2023. [Online]. Available: <https://arxiv.org/abs/2311.08734>
- [55] H. S. Zheng, S. Mishra, X. Chen, H.-T. Cheng, E. H. Chi, Q. V. Le, and D. Zhou, "Take a step back: Evoking reasoning via abstraction in large language models," 2024. [Online]. Available: <https://arxiv.org/abs/2310.06117>
- [56] C. Li, J. Wang, Y. Zhang, K. Zhu, W. Hou, J. Lian, F. Luo, Q. Yang, and X. Xie, "Large language models understand and can be enhanced by emotional stimuli," 2023. [Online]. Available: <https://arxiv.org/abs/2307.11760>
- [57] A. Lu, H. Zhang, Y. Zhang, X. Wang, and D. Yang, "Bounding the capabilities of large language models in open text generation with prompt constraints," 2023. [Online]. Available: <https://arxiv.org/abs/2302.09185>
- [58] Y. Deng, W. Zhang, Z. Chen, and Q. Gu, "Rephrase and respond: Let large language models ask better questions for themselves," 2024. [Online]. Available: <https://arxiv.org/abs/2311.04205>
- [59] Z. M. Wang and et al., "Rolellm: Benchmarking, eliciting, and enhancing role-playing abilities of large language models," 2024. [Online]. Available: <https://arxiv.org/abs/2310.00746>
- [60] M. Yasunaga, X. Chen, Y. Li, P. Pasupat, J. Leskovec, P. Liang, E. H. Chi, and D. Zhou, "Large language models as analogical reasoners," 2024. [Online]. Available: <https://arxiv.org/abs/2310.01714>
- [61] D. Schreiter, "Prompt engineering: How prompt vocabulary affects domain knowledge," Ph.D. dissertation, Georg-August-Universität Göttingen, 2024.
- [62] OpenAI, "Chatgpt: Optimizing language models for dialogue," 2023. [Online]. Available: <https://openai.com/blog/chatgpt>
- [63] W. Alsaweed and S. Aljebreen, "Investigating the accuracy of chatgpt as a writing error correction tool," *International Journal of Computer-Assisted Language Learning and Teaching*, vol. 14, pp. 1–18, 12 2024.
- [64] T. Fang, S. Yang, K. Lan, D. F. Wong, J. Hu, L. S. Chao, and Y. Zhang, "Is chatgpt a highly fluent grammatical error correction system? a comprehensive evaluation," 2023. [Online]. Available: <https://arxiv.org/abs/2304.01746>
- [65] J. Liu, C. S. Xia, Y. Wang, and L. Zhang, "Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation," in *Proceedings of the 37th International Conference on Neural Information Processing Systems*, ser. NIPS '23. Red Hook, NY, USA: Curran Associates Inc., 2023.
- [66] A. Liu and et al., "Deepseek-v3 technical report," 2025. [Online]. Available: <https://arxiv.org/abs/2412.19437>
- [67] B. H. et al., "Qwen2.5-coder technical report," 2024. [Online]. Available: <https://arxiv.org/abs/2409.12186>
- [68] Meta AI, "Llama 3.3 70b instruct model card," <https://huggingface.co/meta-llama/Llama-3.3-70B-Instruct>, 2024, accessed: 2025-05-30.
- [69] OpenAI, "Openai o3-mini system card," <https://openai.com/index/o3-mini-system-card/>, 2025, accessed: 2025-05-30.
- [70] together.ai, "Together models," 2025. [Online]. Available: <https://www.together.ai/models>
- [71] M. A. Covington and J. D. McFall, "Cutting the gordian knot: The moving-average type-token ratio (mattr)," *Journal of Quantitative Linguistics*, vol. 17, no. 2, pp. 94–100, 2010.
- [72] J. P. Kincaid, R. P. Fishburne, R. L. Rogers, and B. S. Chissom, "Derivation of new readability formulas (automated readability index, fog count and flesch reading ease formula) for navy enlisted personnel," Naval Technical Training Command Millington TN Research Branch, Tech. Rep. 8-75, 1975.
- [73] R. Gunning, "The technique of clear writing," (*No Title*), 1952.
- [74] R. Flesch, "A new readability yardstick," *Journal of applied psychology*, vol. 32, no. 3, p. 221, 1948.
- [75] F. Lehner, "Quality control in software documentation: Measurement of text comprehensibility," *Information & Management*, vol. 25, no. 3, pp. 133–146, 1993. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/0378720693900365>
- [76] C. Spearman, "The proof and measurement of association between two things," *The American Journal of Psychology*, vol. 15, no. 1, pp. 72–101, 1904.
- [77] D. J. Hilton, "Logic and causal attribution," in *Part of this chapter is the text of a paper read at the symposium, Attitudes and Attribution: A Symposium in Honour of Jos Jaspars, convened at the Annual Conference of the British Psychological Society, Sheffield, England, Apr 3-5, 1986*. New York University Press, 1988.
- [78] A. Jacovi and Y. Goldberg, "Towards faithfully interpretable nlp systems: How should we define and evaluate faithfulness?" *arXiv preprint arXiv:2004.03685*, 2020.
- [79] S. Gao, X.-C. Wen, C. Gao, W. Wang, H. Zhang, and M. R. Lyu, "What makes good in-context demonstrations for code intelligence tasks with llms?" in *2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 2023, pp. 761–773.
- [80] S. Goriparthi, "Advancing conversational ai: Best practices in prompt engineering for enhanced chatbot performance," *Journal ID*, vol. 6202, p. 8020.
- [81] T. Zimmermann, "Card-sorting: From text to themes," in *Perspectives on data science for software engineering*. Elsevier, 2016, pp. 137–141.
- [82] G. L. Iverson, *Z Scores*. New York, NY: Springer New York, 2011, pp. 2739–2740. [Online]. Available: https://doi.org/10.1007/978-0-387-79948-3_1263
- [83] B. G. Glaser and Hon., "Open coding descriptions," *Grounded Theory Review: An International Journal*, vol. 15, 2016. [Online]. Available: <https://api.semanticscholar.org/CorpusID:149583839>
- [84] "GitHub - prompt-study/prompt-tasks-study — github.com," <https://github.com/prompt-study/prompt-tasks-study/tree/main>, [Accessed 31-05-2025].