

ChatGPT for Code Refactoring: Analyzing Topics, Interaction, and Effective Prompts

Eman Abdullah AlOmar*, Luo Xu*, Sofia Martinez*, Anthony Peruma†, Mohamed Wiem Mkaouer‡, Christian D. Newman§, Ali Ouni¶

*Stevens Institute of Technology, Hoboken, NJ, USA

†University of Hawaii at Manoa, Honolulu, HI, USA

‡University of Michigan-Flint, Flint, MI, USA

§Rochester Institute of Technology, Rochester, NY, USA

¶ETS Montreal, University of Quebec, Montreal, QC, Canada

{ealomar, lxu41, smartine1}@stevens.edu, peruma@hawaii.edu, mmkaouer@umich.edu, cdnvse@rit.edu, ali.ouni@etsmtl.ca

Abstract—Large Language Models (LLMs), such as ChatGPT, have become widely popular and widely used in various software engineering tasks such as refactoring, testing, code review, and program comprehension. Although recent studies have examined the effectiveness of LLMs in recommending and suggesting refactoring, there is a limited understanding of how developers express their refactoring needs when interacting with ChatGPT. In this paper, our goal is to explore interactions related to refactoring between developers and ChatGPT to better understand how developers identify areas for improvement in code, and how ChatGPT addresses developers’ needs. Our approach involves text mining 715 refactoring-related interactions from 29,778 ChatGPT prompts and responses, as well as the analysis of developers’ explicit refactoring intentions. Our results reveal that (1) refactoring interactions between developers and ChatGPT encompass 25 themes including ‘Quality’, ‘Objective’, ‘Testing’, and ‘Design’, (2) ChatGPT’s use of affirmation phrases such as ‘certainly’ regarding refactoring decisions, and apology phrases such as ‘apologize’ when resolving refactoring challenges, and (3) our refactoring prompt template enables developers to obtain concise, accurate, and satisfactory responses with minimal interactions. We envision our results enhancing researchers and practitioners understanding of how developers interact with LLMs during code refactoring.

I. INTRODUCTION

The recent advances in Artificial Intelligence (AI) are revolutionizing computing education in general and software engineering in particular. In fact, the ability of Large Language Models (LLMs) to harness massive amounts of multimodal information has enabled them to perform a variety of tasks that are known to depend on human intervention [1], [2]. Leveraging information from open-source software (OSS) repositories, these LLMs are emerging as assistive technologies for developers in various software engineering tasks, such as code search [3], code quality [4], repair [5], program comprehension [6], generation [7], completion [8], and translation [9]. The promising results of LLMs, in general, and ChatGPT, in particular, have grown in popularity quickly within the software engineering community. A recent GitHub survey with 500 US-based developers [10], shows that up to 92% have AI support integrated into their development environments, and 70% reported an improvement in their coding productivity. Similarly, recent research highlights the

outstanding performance of these models when tested against traditional solutions using existing benchmarks [11].

Despite the existence of built-in models in Integrated Development Environments (IDEs), many developers use conversational models, which are designed to receive text-based requests (*i.e.*, prompts) and generate human-like output texts. The interactive nature of such models has increased their popularity among developers, as they can respond to a broader range of queries besides just code. Built-in models, such as GitHub Copilot¹, provide laser-focused recommendations on current coding tasks, but conversational models, such as ChatGPT, can assist with a wider range of tasks. Therefore, recent studies have shifted from evaluating the capability of ChatGPT, to analyzing collaborative practice and prompt engineering, to improve how developers can foster adequate information that matches their expectations [13], [14].

While LLMs have been heavily solicited for a variety of tasks, little is known about their response to prompts related to code refactoring. Refactoring, by definition, involves enhancing the internal structure of code without changing its external behavior [15], [16]. Due to the subjective nature of refactoring, where several equally valid solutions might exist for a given scenario [15], it is intriguing to observe how language models handle refactoring requests and which quality attributes they prioritize during code optimization. Although there are emerging studies analyzing the refactoring capability of ChatGPT [17], [18], they are mainly focused on benchmarking ChatGPT’s ability to correctly recommend refactorings. Consequently, it is essential to conduct a more detailed examination to comprehend developers’ expectations regarding AI-aided refactoring. Additionally, prior research has not explored situations where developers openly show disappointment with the model’s output.

Figure 1 presents an illustration regarding an issue filed to revise the `ingest()` method. In response to a prompt, ChatGPT delivered an updated code version, clearly stating that the refactoring aimed to enhance *maintainability*. This purpose is also reflected in the concluding title of the pull

¹<https://github.com/features/copilot>

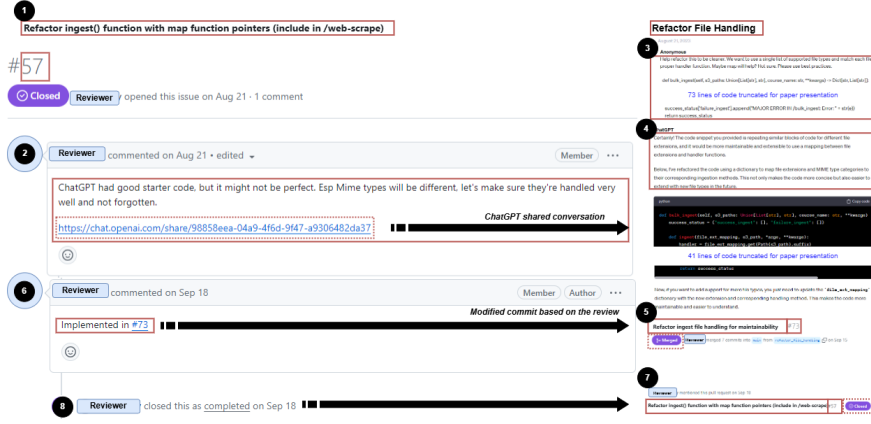


Fig. 1: An illustration of how a developer interacts with ChatGPT and how the generated outcomes have been included in a pull request that was accepted and deployed in production [12].

request that eventually integrated the changes into the live codebase. From this instance, it is evident that ChatGPT targets *maintainability* as a quality attribute in its code optimization process, with the model confirming its quality improvements by using the word ‘certainly.’ Thus, we are motivated to further extract and analyze additional quality attributes to gain a deeper understanding of the refactoring strategies utilized by ChatGPT.

The goal of this paper is to explore what developers care about when they request ChatGPT to refactor their code. Specifically, we aim to identify, categorize, and analyze the patterns in developers’ prompts, to extract the main criteria developers rely on to develop a decision about accepting or rejecting the model’s response. Furthermore, we synthesize this extracted knowledge to craft a prompt template that can accurately reflect developers’ expectations. Our study is based on analyzing 715 refactoring-related interactions collected from 29,778 publicly shared conversations between developers and ChatGPT. These conversations can be found attached to their posted GitHub commits, issues, and files using ChatGPT sharable links. The existence of such traces of interactions is evidence of how developers are trending toward using model output to address their problems, and for possible code alternatives. According to our goal, we aim to answer the following research questions:

- **RQ₁:** *What conversation topics ChatGPT takes into consideration when generating responses to refactoring requests?* To answer this research question, we aim to extract all the coding contexts in which refactoring was performed, whether directly requested by the developer or proposed by ChatGPT as a response. To do so, we conducted a thematic analysis of 715 interactions to design a taxonomy of refactoring contexts. The taxonomy contains 4 root themes, with a total of 25 sub-themes.
- **RQ₂:** *In which scenarios does ChatGPT provide unsatisfactory responses during the refactoring conversations?* In this research question, we are interested in exploring which scenarios ChatGPT’s response was found to be confusing or unsatisfactory. Therefore, we analyze the

model’s responses that tend to include an apology or an affirmation statement during the refactoring interactions.

- **RQ₃:** *What is the performance of the proposed prompt, compared to the dataset?* We address this research question by designing a prompt template that would efficiently refactor the input code. We evaluate the efficiency of the template by assessing its ability to reproduce in *one shot* the exact refactored code that took developers several rounds of prompts to reach it. We also measure our prompt’s length and compare it with the length of the conversations to see if it can lead to the same result with fewer tokens.

Our dataset and artifacts are available for replication and extension purposes [19].

II. STUDY DESIGN

A. Data Collection and Curation

This research employs the DevGPT dataset [20], which comprises extensive data on open-source projects, including code files, commits, issues, Hacker News, pull requests, and discussions. Excluding Hacker News, these elements are all hosted and integrated on GitHub. To collect data from multiple sources, we utilize the following approach:

Step #1: Data Collection: Initially, the data was collected from the DevGPT dataset [20]. This dataset is composed of various JSON files structured into snapshots.

Step #2: Data Extraction: JSON files were retrieved and organized into distinct categories based on their source type. A summary of the extracted data is as follows: Shared ChatGPT links (4,733), GitHub or Hacker News references (3,559), ChatGPT prompts and response (29,778), Code snippets (19,106), Refactoring GitHub commits (470), Refactoring GitHub issues (69), and Refactoring GitHub code files (176).

Step #3: Data Transformation: The JSON files underwent processing to be converted into a structured, relational table format before being imported into the database.

Step #4: Data Translation: Given that some interactions were in non-English languages, we used the Google Translate library to assist in converting non-English content into

English. Such multilingual analysis enhances the clarity and accessibility of the dataset. Once the translation is complete, the database is updated, in a methodical manner, with the new translated text.

Step #5: Data Preprocessing: Titles, bodies, prompts, and responses from different sources underwent cleaning, had stop-words removed, and were tokenized.

We developed a pipeline to perform these tasks, starting with the JSON files as input and producing the desired subset as a database output. Our pipeline integrates multiple technologies, including *SQLite* for handling data, *FastText* to identify conversations in languages other than English, the *Google Translator* library for translating these dialogues, and *NLTK & Spacy* for detailed cleaning and tokenization. Additionally, *Dask* is used to enhance efficiency through concurrent processing during both cleaning and tokenization stages.

B. Taxonomy Building and Refinement

As our research focuses on refactorings, our examination is confined to sources where refactorings were discussed in developer-ChatGPT exchanges. Initially, we extracted all conversations from the original dataset. To confirm that the software artifacts concerned refactoring, we concentrated on prompts illustrating developers’ intent to apply refactoring (*i.e.*, containing the keyword ‘*refactor*’). Searching for the term ‘*refactor*’ is a common method for identifying refactorings in natural language text, and is used by numerous related works to identify refactoring-related commits and text [21]–[26]. This process culminated in selecting three types of sources: commits, issues, and files. Ultimately, to minimize false positives, our analysis was restricted to the presence of ‘*refactor*’ in the prompt for each source type. We found a few data points that were duplicates or contained invalid links; we excluded all such instances. The initial processing stage removed 13 duplicated conversations from the dataset, alongside 53 invalid ChatGPT links. The procedure led to the examination of 470 commits, 69 issues, and 176 files.

The purpose of the manual analysis is to classify the discussion in developer-ChatGPT interactions. Due to the absence of established taxonomies for tasks associated with refactoring through ChatGPT, we utilized a thematic analysis approach following the guidelines of Cruzes *et al.* [27] when examining ChatGPT responses. Thematic analysis is widely regarded in the Software Engineering literature (*e.g.*, [28], [29]) as a method for detecting and documenting patterns (or “themes”) within a set of descriptive labels, referred to as “codes”. For each refactoring response from ChatGPT, the analysis was conducted through the following steps: i) Initial reading of the responses; ii) Creating preliminary codes (*i.e.*, labels) for each response; iii) Translating codes into themes, sub-themes, and more comprehensive themes; iv) Refining themes to identify possibilities for consolidation; v) Defining and naming the final themes while developing a model of higher-order themes supported by evidence.

The procedures outlined earlier were carried out separately by two of the authors. One author independently labeled the

ChatGPT responses, while the other focused on examining the taxonomy draft. After each cycle, both authors convened to refine the taxonomy. For the manual coding process, we employed a spreadsheet tool with tagging functions. This spreadsheet offered the annotators the following details: (1) the prompt given to ChatGPT, (2) ChatGPT’s response, and (3) a link for ChatGPT file sharing. During the study, one author possessed 7 years of research experience in refactoring, whereas the other author had 11 years of experience in the same field.

It is important to note that the approach is not a single-step process, with each round focusing on approximately one-third of the inspection instances. Once the two authors have inspected all 715 instances, we solved the conflicts in 6% of the cases. Conflicts can be attributed to two design choices. Firstly, the absence of predefined categories was a key factor. This meant that discrepancies arose when two authors independently assigned semantically equivalent, yet distinct, labels to describe a given automated task. Secondly, our cautious approach to defining conflicts contributed to this observation. We define an instance as a conflict if two authors assign different sets of labels, even if there is partial overlap between the two sets. Throughout the analysis of the codes, many first-cycle codes were integrated into other codes, renamed, or completely eliminated. During the process when the authors worked to convert the information into themes, they had to shuffle, fine-tune, and occasionally reclassify data into either different or new codes. For instance, we consolidated the initial categories such as “*incorrect refactoring*”, “*behavior preservation violation*”, “*separation of other changes from refactoring*”, and “*interleaving other changes with refactoring*” into the overarching category of “*Objective*”. This method resulted in the development of 25 themes, which were then used to create a hierarchical taxonomy that covers refactoring tasks found in ChatGPT responses.

C. Refactoring Prompt Template Construction and Evaluation

Step #1: Prompt Construction. In interactions with ChatGPT, developers articulate software development concepts using natural language. Due to the varied methods developers use to describe issues, it is impractical to rely solely on automated methods for analyzing prompts and responses. Thus, we conducted a manual examination to construct a prompt for refactoring in three rounds. Two authors independently reviewed the developer-ChatGPT interactions (*i.e.*, prompt, response, and code of each interactions), constructed and tested the prompt, resolving any discrepancies through discussion. In the first round, two authors equally divided the ChatGPT-developer conversations. Then, we started by manually reviewing the unique set of prompts in which developers initiated a conversation with ChatGPT about refactoring. Based on our observations, we recognized that there is a certain format that leads to a shorter interaction between ChatGPT and developers. The effective format usually involves a precise breakdown of prompt explain to ChatGPT what role it is playing, the working set that the developer is using,

the refactoring task, the project specifications, requirements, output format example, the operating system being used, the installed tools, and how it is being run. Consequently, we construct our first prompt. Then, in the second round, we review the literature on the realm of LLM refactoring [14], [17], [30]–[44] to better understand what prompts have been considered successful. We found that this is a multi-faceted topic, with research expanding into different LLMs, different tools, different prompting techniques, and different refactoring types. From these studies, we noticed that there are certain prompting techniques that are better for refactoring. Thus, we implemented these techniques into the prompt by refining our drafted prompt to include any missing attributes accordingly in our third iteration. The authors evenly split the instances for review. If there were discrepancies, the authors would engage in extensive discussions until a consensus was reached. These discussions often revolved around whether there is noteworthy information to be added in our prompt. Ultimately, both authors agreed on the constructed prompt after discussing each instance.

Step #2: Prompt Evaluation. After constructing the prompt, we evaluated it by reproducing the corresponding developer’s prompt from DevGPT, which serves as our ground truth. The evaluation was conducted using three key metrics: number of turns, prompt length, and response length. To analyze the data, we first tested normality using the Shapiro-Wilk test and found that the distribution of the metrics did not follow a normal distribution. As a result, we applied the Mann-Whitney U test [45], a non-parametric statistical test, to compare the two independent groups. The null hypothesis is defined by the absence of variation in the metric values of the developer’s prompt and our prompt. Thus, the alternative hypothesis indicates that there is a variation in the metric values. Furthermore, the variation between values of both sets is considered significant if its associated p -value is less than 0.05. Furthermore, we use the Cliff’s Delta (δ) [46], a non-parametric effect size measure, to estimate the magnitude of the differences between DevGPT developer prompt and our prompt. Regarding its interpretation, we follow the guidelines reported by Romano *et al.* [47]: Negligible for $|\delta| < 0.147$, Small for $0.147 \leq |\delta| < 0.33$, Medium for $0.33 \leq |\delta| < 0.474$, and Large for $|\delta| \geq 0.474$.

III. EXPERIMENTAL RESULTS

A. RQ1: What conversation topics ChatGPT takes into consideration when generating responses to refactoring requests?

After examining the refactoring responses provided by ChatGPT, we establish broad high-level classifications for the developer-ChatGPT refactoring dialogue. Figure 2 shows the taxonomy resulting from our analysis. The taxonomy is structured into two tiers: the upper tier includes four categories that organize activities with related objectives, while the lower tier comprises 25 subcategories, offering detailed classification. Conversations regarding developer-ChatGPT refactoring revolve around four primary categories illustrated in the figure: (1) quality, (2) objective, (3) testing, and (4) design. It is

important to highlight that our categorization is not exclusive; thus, a response might belong to multiple categories. Examples for each category are listed in Table I. The ensuing part of this subsection delves deeper into these categories.

Category #1: Quality. When addressing refactoring requests, the quality of design plays a crucial role. Based on ChatGPT feedback, it ensures compliance with *coding conventions*, promotes *code reviews*, and maintains *documentation* standards. Additionally, it aims to enhance both *internal* and *external quality attributes*, while steering clear of design pitfalls like *code smells*. For example, ChatGPT advises on the best practices for code writing and the optimization of *internal and external quality attributes*, as developers may not perceive the complete overview of the software design.

Category #2: Objective. This category compiles responses that emphasize assessing the accuracy of code transformations and determining if the proposed modifications result in a secure and reliable refactoring. These ChatGPT responses investigate topics such as *refactoring correctness*, *behavior preservation*, and the appearance of refactoring alongside other changes, in addition to responses related to refactoring operations such as *method composition*, *feature move*, and *generalization*. Since developers frequently interleave refactoring with other activities, ChatGPT highlighted that combining refactoring with additional changes could obscure errors, thus increasing the risk of introducing bugs. Further, we have compiled instances where ChatGPT ultimately requests clear documentation of *feature-related*, *bug fix-related*, and *clean up* tasks to enhance understanding of the rationale behind submitted refactoring requests. This demonstrates ChatGPT’s recurring suggestions for improving developers’ documentation practices, particularly concerning the perception and reasoning of the changes.

Category #3: Testing. The goal of refactoring is to enhance internal software structure while keeping its functionality intact. Ideally, pre-existing unit tests should be adequate to confirm this consistency. However, since refactoring is often mixed with other tasks, alterations in software behavior might occur. In these scenarios, unit tests might fail to notice these changes unless they are updated to accommodate the new features. This issue has been highlighted in several ChatGPT responses, especially when developers overlook these behavior changes. Our analysis of these responses shows that ChatGPT suggests implementing unit tests prior to refactoring to ensure the code remains intact. It also recommends following the red-green-refactor methodology. Additionally, when developers seek guidance on testing during refactoring, ChatGPT recommends updating the test files. Thus, within the *Testing* category, we identify the following sub-categories: *test suite*, *red-green-refactor*, and *production synchronization*.

Category #4: Design. According to ChatGPT responses, developers are asked to follow general guidelines, design pattern, and specific design principles (*e.g.*, SOLID and GRASP). ChatGPT’s responses about the SOLID principle focus on key design principles such as the single responsibility principle, the Liskov substitution principle, and the dependency inversion

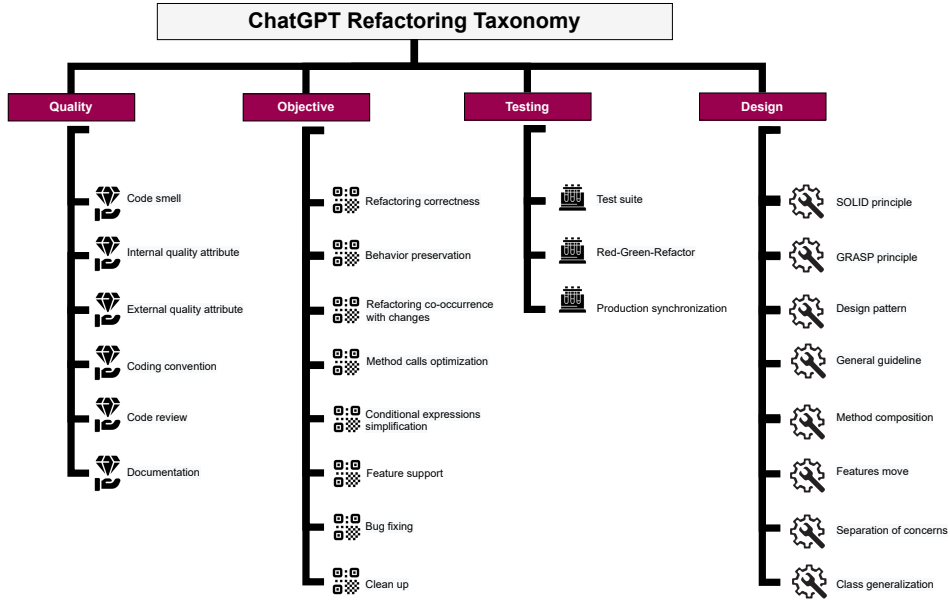


Fig. 2: A taxonomy of the ChatGPT refactoring topics.

principle. On the other hand, ChatGPT’s responses about the GRASP principles offer guidelines for assigning responsibilities to classes such as low-coupling and high-cohesion. Additionally, it discussed refactoring changes such as *method composition*, *feature move*, and *class generalization*. This indicates that ChatGPT asked developers to strive to design classes that are focused and easy to maintain by breaking down complex systems into more manageable components.

B. RQ₂: In which scenarios does ChatGPT provide unsatisfactory responses during the refactoring conversations?

In the refactoring process, ChatGPT’s expressions of regret (such as admitting errors) and confirmations (such as maintaining confidence or agreement) are important. Regrets may be offered if the refactoring leads to errors or unexpected results. Affirmations serve to recognize refactoring successes. Upon analyzing ChatGPT’s responses, we cluster the responses into categories.

Category #1: Apology. When responding to errors in its refactoring suggestions, ChatGPT often employs phrases like ‘My apologies’ and ‘I apologize.’ These expressions serve to acknowledge mistakes and reflect a promise to rectify them. ChatGPT expresses regret when it introduces bugs in the code, misinterprets the code’s content, overlooks update capabilities, incorporates poor practices, or displays misunderstanding in refactoring. ChatGPT response exemplifies this category is shown in Figure 3. As can be seen, ChatGPT apologizes for introducing duplicate code, which is considered one of the bad practices that violates the best design principles. This demonstrates ChatGPT’s awareness of the importance of maintaining code quality and adhering to established coding standards. Next, we provide an analysis of these categories:

Example of apology [48]:

You’re correct! My apologies for the redundancy in the code. We can indeed refactor the ‘**TasksList**’ component to use the ‘**PromptDescriptorViewer**’ component instead of re-implementing its functionality. This would allow us to uphold the DRY (Don’t Repeat Yourself) principle in our code.

Example of affirmation [49]:

Certainly, it’s a good idea to create a separate function for fetching the descriptor URL. We can place this function in a common utilities or services file. Let’s refactor the code to make this change. We will create a new file ‘**apiServices.js**’ for such service functions.

Fig. 3: Chatbot interaction in the context of refactoring.

- **Bugs introduction.** ChatGPT apologizes when its refactoring suggestions introduce bugs or errors into the code. For example, it may suggest a change that inadvertently breaks the functionality of the program. This type of apology is crucial as it reflects the model’s recognition of its mistake and its impact on the code’s functionality.
- **Lack of understanding.** ChatGPT sometimes misinterprets the developer’s code or intent. It apologizes for any misunderstandings about refactoring and tries to clarify or correct its previous suggestions.
- **Missing update functionality.** When ChatGPT fails to account for necessary updates across different parts of the code, it issues an apology. This demonstrates the need for comprehensive changes that span multiple areas of the

TABLE I: Examples of the ChatGPT refactoring topics.

Category	Sub-category	Example (Excerpts from a related refactoring conversations)
Quality	Code smell	"This version of the shell script refactors the existing solution to remove the duplication of the descriptor filename."
	Internal quality attribute	"Currently, many classes are tightly coupled with their dependencies. Try to inject dependencies instead of creating them inside the classes. This would make your classes easier to test and maintain."
	External quality attribute	"Your TopNav.js component is quite modular and you've done a good job organizing the related components. Below, I'll offer some refinements to increase the maintainability and readability of the code."
	Coding convention	"Consistent Naming Conventions: Stick to one naming convention throughout your project for variables, methods, classes etc."
	Code review	"Review your refactored code and validate your use of inheritance."
Objective	Documentation	"when creating a git commit the message is crucial for readability and understanding the purpose of the commit for future reference heres how you could craft a commit message for this refactoring."
	Refactoring correctness	"it's important to be mindful of Liskov's Substitution Principle (LSP), which states that objects of a superclass shall be able to be replaced with objects of a subclass without affecting the correctness of the program."
	Behavior preservation	"The EOF marker is quoted to prevent parameter substitution within the heredoc, preserving the exact content we want for the App.jsx file."
	Refactoring co-occurrences with changes	"For brevity, I'm going to write a simple class demonstrating how you might implement some of these features using Google's JavaScript client library. Keep in mind this is just a starting point and you would need to expand upon this with proper error handling, user interface components, and more."
	Method calls optimization	"Thank you for pointing out those mistakes i apologize for the confusion youre correct if the project is already configured to support es the r esm flag should not be needed and the temp file should be cleaned up heres the corrected scriptcodeblockthis corrected script should take care of the renaming refactoring and the update to packagejson without any lingering temp files."
	Conditional expressions simplification	"Let's refactor the Heading component to be a functional component. This refactor will also simplify the conditional logic and improve readability."
	Feature support	"You'd like to refactor your application's image caching system to use Google Drive instead of your backend server. We'll need to interface with Google Drive's API, and we will need to implement features such as saving image payloads, image files, pagination, pre-signed URLs, and deletion functionality."
	Bug fixing	"You are throwing exceptions but not catching them. Make sure to handle exceptions properly."
	Clean up	"The task involves refactoring and reorganizing the existing code into separate modules and folders for better manageability and separation of concerns."
	Test suite	"Testing: You're using Cypress, Enzyme, and Jest, which is great for testing. Make sure the tests are updated as you refactor."
Testing	Red-Green-Refactor	"Red-Green Refactor is an excellent approach here. Start with a failing test (Red), implement your function, make the test pass (Green), and then refactor if necessary."
	Production synchronization	"Given the refactoring of the Game.js file into a new User.js file, you'll need to refactor the unit tests accordingly."
Design	SOLID principle	"Some classes like ContractRpcWriter do a lot of work - initialization of Web3, gas estimation, transaction signing, and sending. Try to keep a class to a single responsibility. This will increase the modularity and readability of your code."
	GRASP principle	"This design enables you to use polymorphism, one of the core principles in OOP, where a class has many ("poly") forms ("morphs")."
	Design pattern	"You're currently using the Strategy pattern in your design, where the Algorithm class defines a common interface for all strategies (algorithms), and Minimax and UserInput are concrete strategies that implement this interface. This pattern provides a good solution for your current requirements."
	General guideline	"Your refactored code looks much cleaner and modular! Dividing responsibilities between different classes is a good software engineering practice. It simplifies the understanding of code and makes maintenance easier."
	Method composition	"This script will refactor the code by splitting it into separate files. It creates a new file setupRoutes.js to handle all the routing logic, deletes the handlers.js file as it's no longer needed, and moves the generateHandler function to a new file generateHandler.js."
	Features move	"Creates a new directory named backend inside the src directory. Moves server.js and servePromptDescriptor.js to the backend directory. Refactors server.js by extracting the route handlers into a new file called handlers.js."
	Separation of concerns	"The user class now encapsulates all the user-related functionalities. getStoredUser, getPicks, updatePicks, backfillResults, switchUser, and displayBackfilledResults functions are now methods of the User class. This provides a clear separation of user-specific logic and game logic."
	Class generalization	"Review your refactored code and validate your use of inheritance. Make sure that it adheres to Liskov's Substitution Principle (LSP). LSP is a concept in object-oriented programming that states that if a program is using a base class, it should be able to use any of its subclasses without the program knowing it."

codebase.

- **Bad practice introduction.** ChatGPT's suggestions may violate coding best practices. In such cases, it acknowledges the mistake and attempts to rectify it.

Category #2: Affirmation. In its responses, ChatGPT includes keywords like 'certainly', 'absolutely', 'of course', and 'sure' to affirm the developer's prompt on refactoring. These words reflect a high level of agreement or assurance in addressing the developer's question or suggestion. ChatGPT consistently expresses this degree of confidence or agreement throughout its refactoring suggestions. When recommending refactoring, ChatGPT aims to improve code quality, offers examples for code refactoring, adheres to established refactoring standards or principles, and follows best coding practices, such as modularization and abstraction, to boost code maintainability and readability. Additionally, it introduces alternative refactoring methods to give developers a range of strategies or techniques to explore. Figure 3 presents an example of the ChatGPT responses that exemplify this category. As shown, ChatGPT adheres to the best practice of creating separate functions with common logic. Below, we provide an analysis of these categories:

- **Code quality improvement.** ChatGPT affirms its goal of improving code quality through its suggestions. This includes applying best coding practices to enhance maintainability and readability.
- **Refactoring demonstration.** When affirming a refactoring suggestion, ChatGPT provides concrete examples and guidance, demonstrating how to implement the refactoring effectively.
- **Code guideline adherence.** ChatGPT affirms its understanding of the developer's intent by providing guidelines. This reassurance helps in building trust and ensures that the developer feels understood.
- **Best practice application.** ChatGPT often affirms its adherence to best practices in its suggestions. This is important for ensuring that the code remains maintainable and follows coding standards.
- **Refactoring extension.** In addition to affirming a single refactoring suggestion, ChatGPT sometimes suggest applying additional refactorings.

In the course of this analysis, we identified several potential problems with the responses, which were not directly related to the refactoring request. In fact, we witnessed several instances

of issues encompassing reliance on outdated information, technical errors, incorrect API calls, and fabricated source classes. Although these mistakes can have significant impacts, we are not evaluating their severity in this study, as our focus is on the refactoring context. Investigating the severity of these errors is an intriguing subject that would necessitate more information, such as the complete source code of the projects. However, such data is not available in the current dataset we utilize. So, as part of future work, we plan to gather such relevant information and conduct a severity analysis since it would be interesting to investigate whether refactoring would increase the probability of these mistakes.

C. RQ₃: What is the performance of the proposed prompt, compared to the dataset?

Prompt Engineering is a technique for guiding LLMs to generate specific outputs. The need for studying prompting in LLM-driven code refactoring comes from the need to improve its effectiveness, and reliability. While LLMs have shown promise in automating refactoring tasks [30], their performance is highly dependent on how prompts are structured and syntactically formed. Through studies, different refactoring tasks require varying levels of context and reasoning. Simpler refactorings like renaming can often be handled with basic prompts. However, more complex refactorings, such as the extractions, demand more structure, iterative, and context-based prompts to ensure correctness.

Our goal is to guide developers by creating refactoring prompt template that can reduce time and effort, allowing the model to adapt to a variety of refactoring tasks. Figure 4 illustrates our proposed prompt template for conversational refactoring using LLMs. The crafted template consists of the following main sections:

- **Role.** The prompt starts by specifying a role as an AI expert in refactoring and software quality. Reynolds and McDonell mentioned that persona-based pattern (*i.e.*, AI is assigned a specific role), significantly improves contextual understanding [50].
- **Working Set.** This section allows developers to provide their code for refactoring purposes.
- **Context.** This section allows developers to indicate the context pertaining to the language, project, installed tools, and constraints.
- **Refactoring Task.** This section is the main part that allows developers to provide specific refactoring tasks, including motivation and intended impact [23]. Since recent study has shown that developers often make generic refactoring requests, while ChatGPT typically includes the refactoring intention [14], our goal is to assist developers in utilizing the right information in the prompt so that LLMs can focus on the intended improvement instead of making random changes.
- **Steps to Follow.** This section ensures a step by step instruction of how LLMs should proceed with refactoring task, while ensure behaviour preservation and successful test results.

Note for developer: text written in blue is meant for your eyes only, not as prompt for chat. Items in <> should be filled in as your own texts, not to keep.

Prompting techniques: few shots, one shot, zero shot, context constraints, output constraints, chain of thought.

You are an expert in programming, refactoring, and providing advice in software quality. Given the following prompt/-code snippet, modify/write the code so that it best matches what the user wants, and provide comments explaining your work. If needed, ask clarifying questions to check both your and the user's understanding of the prompt.

Working Set:

<here you will input in the code you are working with. Make sure to state clearly if they are from the same file or different files to avoid confusion.>

Context:

- Language: *<e.g., JavaScript, Python, Java>*
- Project: *<brief description of how/where this code is used>*
- Installed Tools: *<list dependencies, libraries, frameworks>*
- Constraints: *<e.g., Do not change the output behavior; Maintain API structure, output constraints>*

Refactoring Task:

<give refactoring prompt as much detail as possible, give two - three full sentences MINIMUM. Use these keywords as a starting context, these are meant to be used in the beginning of the action:

Refactoring Patterns (Self-Affirmed Refactoring)

<refactoring pattern-related quality issues words such as these, after the initial action. Start with "to improve...":>

Refactoring Intent (Self-Affirmed Refactoring)

- <explicitly state any components that you want/would prefer chat to use.>

<To improve the code:>

- <Clearly state any components or patterns that should be used.>

- <Ensure modularization where applicable.>

- <Follow established design patterns.>

Steps to Follow:

1. Analyze the Code: Identify areas that require refactoring.
2. Apply Refactoring Strategies: Improve quality attributes (e.g., readability, maintainability, efficiency, performance, cohesion, coupling, etc).
3. Provide Multiple Solutions: Generate at least two refactored versions with explanations.
4. Validate Functionality: Ensure the new code behaves the same as the original.
5. Comment & Document: Provide commit messages and inline comments explaining the changes.

Output Format:

- Provide refactored code inside a code block.
- Include a before-and-after comparison.
- Provide concise commit messages summarizing each major change.
- Generate test cases to validate the refactoring.

Example start:

<how it is being ran>

Example end

If there are any unclear instructions, please ask clarifying questions.

Fig. 4: Refactoring prompt template.

TABLE II: Statistics of LLM prompt efforts.

Metrics	Developer's prompt						Our prompt						Statistical difference	
	Min	Q1	Median	Mean	Q3	Max	Min	Q1	Median	Mean	Q3	Max	p-value	Cliff's delta (δ)
Number of turns	1	2	3	13.58	9.50	11	1	1	1	1.45	2	3	0.000005649	large (0.56)
Prompt length	372	2996	4017	16604	7812	13537	1563	3530	4489	4855.52	5287.50	6570	0.1952	small (0.04)
Response length	1897	3865	5871	27538.23	21049.50	23043	2297	4192	5778	6975.71	7630	11445	0.6376	small (0.11)

- **Output Format.** This section specifies the output format of the refactored code requested by the developer.
- **Example.** This section specifies the example provided by the developer.
- **Clarification Option.** This part ensures LLMs understanding of the information given in the prompt.

Regarding the evaluation of the proposed prompts, we run our prompt on the same code provided by the original DevGPT developers and calculate the three metrics. Looking at Table II, we found that our prompts differ (*i.e.*, fewer turns ($\mu = 1.45$), shorter prompt ($\mu = 4855.52$), and shorter response ($\mu = 6975.71$)) from original developer's prompts. We also performed a non-parametric Mann-Whitney U test and we obtained a statistically significant p -value when the values of these two groups were compared (p -value < 0.05 for number of turns), and accompanied with a small or large effect size depending on the prompt effort/metric. We notice that usually when the developer provides the information listed in our prompt, the refactoring conversations are shorter and more fruitful.

IV. DISCUSSION

Takeaway #1: Effective prompting of ChatGPT for minimal-interaction refactoring responses remains an area for further exploration. Since our study aims to constitute a “successful” prompt in the context of refactoring, we propose a structured prompt template to help developers get a satisfactory response from GPT with the least possible shots (*i.e.*, prompt rounds). This effective structured prompt pattern serves as a guideline for prompting related to refactoring and may have the potential to be very useful to practitioners who use GPT for refactoring tasks. For example, it can improve accuracy, relevance, developer productivity, and save time and effort (RQ₃).

Takeaway #2: LLM-Driven refactoring is sub-optimal for certain refactoring contexts. Our taxonomy shows that developers prompted LLMs for various refactoring operations, ranging from method composition, feature move, data organization, simplifying method calls, to dealing with generalization. This sheds light on exploring two important directions: (1) what refactoring operations that LLMs can/cannot handle. This was observed in the apologetic and affirmation phrases due to the challenge of LLMs can only maintain a limited understanding of the broader context (RQ₂), and (2) what refactoring operations are considered underrepresented (RQ₁). For apologetic cases, we found that ChatGPT struggled the most with situations involving bug introduction, lack of understanding, and the introduction of bad practices. In the DevGPT example [51], the developer had to prompt GPT 11 times to obtain a satisfactory response, with ChatGPT repeatedly

apologizing for misunderstandings, oversights, and confusion. Consequently, we believe that developers can change their prompting strategy to get a satisfactory result with fewer turns by improving prompts with clarity, requesting GPT to generate multiple solutions, and asking for reasoning/explanation behind the solutions. These strategies have been incorporated into our prompt template, which acts as a guide for developers when performing conversational refactoring.

Takeaway #3: LLM-generated code can be unreliable; requiring strict processes and methods to detect and highlight unreliable, generated code. According to our analysis, LLMs can provide correct refactoring-related responses. However, there are cases where LLMs can hallucinate or suggest sub-optimal solutions, which especially shown in Developer-ChatGPT refactoring conversation (RQ₁). According to RQ₃, while compared to the developer prompts, we had fewer iterations of prompting needed, there were still certain issues that come from automating refactoring using ChatGPT. Issues that have been mentioned by previous study [52], such as ChatGPT-4's abilities in adhering to formatting based off the requirements set by the developers. For example, in multiple of our tests, when it involves ChatGPT to write markdown files, such as README.mds [51], ChatGPT rarely succeeds in doing so. We analyze the failures that came from such behaviors and find that it is likely that ChatGPT is unaware of when and where it has left a markdown code box. This leads it to create awkward instances of responses where information that should be retained in markdown jumps out of the markdown block, and instances where ChatGPT was adding an afterword in its response but ended up in a code block.

Outside of code formatting, we also noticed that in instances where we ask ChatGPT to provide commit messages, it often fails to do so in Git format (*e.g.*, git commit -m “Hello World!”). This could mean various issues, such as ChatGPT not understanding the needs of the developers or having formatting issues again. ChatGPT often fails to follow the entire prompt. One of the common mistakes made by ChatGPT was skipping creating test cases and creating multiple versions of the code and assessing them by comparisons. We had made those two requirements a stable in our prompt formatting, but most of the time when tested ChatGPT either completely ignores it or asks us to prompt again if we want the test cases. Additionally, when it comes to ChatGPT's refactoring suggestions, it has the potential to hallucinate. For instance, when applied to real-world Extract Method instances, ChatGPT-3.5 was found to produce two types of hallucinations: providing refactoring suggestions with either syntactical errors or illogical code that does not address the issue, such as suggesting to extract an entire method body. When reviewing the original developer prompts from DevGPT, we found that ChatGPT sometimes produced illogical code when asked to produce markdown files, and it exhibited another hallucination where it misunderstood the location of a file within a project's file structure, both causing developers to make additional prompts to address these issues. Developers can take steps

to minimize these hallucinations by using tools that filter out incorrect or impractical suggestions. For example, Pomian *et al.* [53] developed EM-ASSIST that uses LLMs to generate, validate, enhance, and rank suggestions. The EM-ASSIST IntelliJ plug-in generates 5-10 suggestions for the selected method, employing program slicing (*e.g.*, excluding code that does not affect a target variable or statement) to enhance the generated suggestions, then using the IntelliJ IDEA API to filter invalid suggestions by verifying whether the code meets the refactoring preconditions for Extract Method. Finally, the plug-in ranks the suggestions based on how frequent a particular refactoring is suggested by the LLM, and presents them to the developer to select from. Thus, developers can use tools such as EM-ASSIST to reduce the amount of suggestions with hallucinations, while also spending less time finding the most useful suggestions for their issue. Ultimately, further research exploring and addressing these types of hallucinations is needed to increase the reliability of ChatGPT’s suggestions and LLM-based refactoring in general.

V. THREATS TO VALIDITY

External Validity. We center our study on open-source systems. Consequently, our findings may not be generalized to all other open-source systems or commercially developed projects. Moreover, although our dataset covers different programming languages, our results may not generalize to systems written in other languages that are not considered in this study. Nevertheless, the objective of this paper is not to construct a theory universally applicable to all systems, but to demonstrate the capability of ChatGPT to provide solutions for code refactoring. Another potential concern pertains to the proposed taxonomy. It may not apply universally to other open-source or commercial projects. Therefore, we cannot claim that the findings regarding ChatGPT refactoring taxonomy are applicable to all software systems, especially those where the imperative for design enhancement might be less significant. Finally, our study is solely focused on ChatGPT, and it is recommended that future research extend to other general-purpose chatbots applicable to software development, such as the recently released Google Bard, which may have limitations that ChatGPT does not possess [54].

Internal and Construct Validity. Concerning the identification of developer-ChatGPT refactoring conversations, we constructed our dataset by extracting data source types that contain the term ‘refactor’ in the prompt. There is the possibility that we have excluded synonymous terms/phrases. However, even though this approach reduces the number of source types in our dataset, it also decreases false positives, and ensures that we analyze artifacts that are explicitly focused on refactorings. Another potential concern regarding validity involves the developer-ChatGPT refactoring conversation. Given that refactorings may be interleaved with other modifications (*i.e.*, developers executed changes alongside refactorings) [55], we cannot assert that the chosen refactoring conversations solely pertain to refactoring activities. Moreover, throughout the manual analysis, we intentionally omitted cases where

ChatGPT’s contribution to refactoring activity was ambiguous, and we applied multiple labels in cases where ChatGPT served multiple purposes. Finally, we designed our refactoring prompt template with generalizability in mind but tested it with GPT to be consistent with DevGPT dataset, its effectiveness still need to be assessed using other LLMs.

VI. CONCLUSIONS

Large Language Models (LLMs) have seen a rise in popularity and are extensively utilized in numerous software engineering tasks. These tasks include, but are not limited to, refactoring, testing, code review, and program comprehension. This paper aims to explore interactions between developers and ChatGPT concerning refactoring, aiming to gain insight into how developers identify areas for code improvement and how ChatGPT addresses their needs. Our approach consists of extracting refactoring-related discussions from 29,778 ChatGPT prompts and responses, with an emphasis on developers’ clear refactoring objectives. Our results show that (1) developer-ChatGPT refactoring interactions cover 25 themes, (2) ChatGPT’s affirmations and apologies in multi-turn conversations highlight failure points, helping us study when and how it falls short, and (3) our refactoring prompt template helps developers receive precise, efficient, and satisfactory responses with the least amount of interaction.

Declaration of generative AI and AI-assisted technologies in the writing process. During the preparation of this work, the author used the ChatGPT Web interface and Writefull to improve the language and readability of some sections. After using this tool, the authors reviewed and edited the content as needed and take full responsibility for the content of the publication.

REFERENCES

- [1] N. Nathalia, A. Paulo, and C. Donald, “Artificial intelligence vs. software engineers: An empirical study on performance and efficiency using chatgpt,” in *Proceedings of the 33rd Annual International Conference on Computer Science and Software Engineering*, pp. 24–33, 2023.
- [2] A. Ahmad, M. Waseem, P. Liang, M. Fahmideh, M. S. Aktar, and T. Mikkonen, “Towards human-bot collaborative software architecting with chatgpt,” in *Proceedings of the 27th International Conference on Evaluation and Assessment in Software Engineering*, pp. 279–285, 2023.
- [3] Y. Wanga, L. G. E. Shic, W. C. J. Chena, W. Z. M. W. H. Lie, and H. Z. Z. L. Z. Zhenga, “You augment me: Exploring chatgpt-based data augmentation for semantic code search,”
- [4] J. White, Q. Fu, S. Hays, M. Sandborn, C. Olea, H. Gilbert, A. El-nashar, J. Spencer-Smith, and D. C. Schmidt, “A prompt pattern catalog to enhance prompt engineering with chatgpt,” *arXiv preprint arXiv:2302.11382*, 2023.
- [5] M. A. Haque and S. Li, “The potential use of chatgpt for debugging and bug fixing,” *EAI Endorsed Transactions on AI and Robotics*, vol. 2, no. 1, pp. e4–e4, 2023.
- [6] W. Ma, S. Liu, W. Wang, Q. Hu, Y. Liu, C. Zhang, L. Nie, and Y. Liu, “The scope of chatgpt in software engineering: A thorough investigation,” *arXiv preprint arXiv:2305.12138*, 2023.
- [7] Y. Feng, S. Vanam, M. Cherukupally, W. Zheng, M. Qiu, and H. Chen, “Investigating code generation performance of chat-gpt with crowd-sourcing social data,” in *Proceedings of the 47th IEEE Computer Software and Applications Conference*, pp. 1–10, 2023.
- [8] R. Pudari and N. A. Ernst, “From copilot to pilot: Towards ai supported software development,” *arXiv preprint arXiv:2303.04142*, 2023.
- [9] W. Jiao, W. Wang, J.-t. Huang, X. Wang, and Z. Tu, “Is chatgpt a good translator? a preliminary study,” *arXiv preprint arXiv:2301.08745*, 2023.

- [10] <https://futurism.com/the-byte/github-92-percent-programmers-using-ai>.
- [11] W. Sun, C. Fang, Y. You, Y. Miao, Y. Liu, Y. Li, G. Deng, S. Huang, Y. Chen, Q. Zhang, *et al.*, "Automatic code summarization via chatgpt: How far are we?," *arXiv preprint arXiv:2305.12865*, 2023.
- [12] <https://github.com/UIUC-Chatbot/ai-ta-backend/issues/57>.
- [13] H. Hao, K. A. Hasan, H. Qin, M. Macedo, Y. Tian, S. H. Ding, and A. E. Hassan, "An empirical study on developers shared conversations with chatgpt in github pull requests and issues," *arXiv preprint arXiv:2403.10468*, 2024.
- [14] E. A. AlOmar, A. Venkatakrishnan, M. W. Mkaouer, C. Newman, and A. Ouni, "How to refactor this code? an exploratory study on developer-chatgpt refactoring conversations," in *Proceedings of the 21st International Conference on Mining Software Repositories*, pp. 202–206, 2024.
- [15] M. Fowler, K. Beck, J. Brant, W. Opdyke, and d. Roberts, *Refactoring: Improving the Design of Existing Code*. Boston, MA, USA: Addison-Wesley Longman Publishing Co., Inc., 1999.
- [16] T. Mens and T. Tourwé, "A survey of software refactoring," *IEEE Transactions on software engineering*, vol. 30, no. 2, pp. 126–139, 2004.
- [17] D. Pomian, A. Bellur, M. Dilhara, Z. Kurbatova, E. Bogomolov, T. Bryksin, and D. Dig, "Next-generation refactoring: Combining llm insights and ide capabilities for extract method," in *2024 IEEE International Conference on Software Maintenance and Evolution (ICSME)*, pp. 275–287, IEEE, 2024.
- [18] B. Liu, Y. Jiang, Y. Zhang, N. Niu, G. Li, and H. Liu, "Exploring the potential of general purpose llms in automated software refactoring: an empirical study," *Automated Software Engineering*, vol. 32, no. 1, p. 26, 2025.
- [19] <https://smilevo.github.io/self-affirmed-refactoring/>.
- [20] T. Xiao, C. Treude, H. Hata, and K. Matsumoto, "Devgpt: Studying developer-chatgpt conversations," in *2024 IEEE/ACM 21st International Conference on Mining Software Repositories (MSR)*, pp. 227–230, IEEE, 2024.
- [21] E. Murphy-Hill, A. P. Black, D. Dig, and C. Parnin, "Gathering refactoring data: a comparison of four methods," in *Proceedings of the 2nd Workshop on Refactoring Tools*, pp. 1–5, 2008.
- [22] E. A. AlOmar, M. W. Mkaouer, and A. Ouni, "Can refactoring be self-affirmed? an exploratory study on how developers document their refactoring activities in commit messages," in *International Workshop on Refactoring-accepted*. IEEE, 2019.
- [23] Z. Di, B. Li, Z. Li, and P. Liang, "A preliminary investigation of self-admitted refactorings in open source software (s)," in *International Conferences on Software Engineering and Knowledge Engineering*, vol. 2018, pp. 165–168, KSI Research Inc. and Knowledge Systems Institute Graduate School, 2018.
- [24] J. Ratzinger, T. Sigmund, and H. C. Gall, "On the relation of refactorings and software defect prediction," in *Proceedings of the 2008 International Working Conference on Mining Software Repositories, MSR '08*, (New York, NY, USA), pp. 35–38, ACM, 2008.
- [25] E. A. AlOmar, M. W. Mkaouer, and A. Ouni, "Toward the automatic classification of self-affirmed refactoring," *Journal of Systems and Software*, vol. 171, p. 110821, 2021.
- [26] E. A. AlOmar, H. AlRubaye, M. W. Mkaouer, A. Ouni, and M. Kessentini, "Refactoring practices in the context of modern code review: An industrial case study at xerox," in *2021 IEEE/ACM 43rd International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP)*, pp. 348–357, IEEE, 2021.
- [27] D. S. Cruzes and T. Dyba, "Recommended steps for thematic synthesis in software engineering," in *2011 international symposium on empirical software engineering and measurement*, pp. 275–284, IEEE, 2011.
- [28] D. Silva, N. Tsantalis, and M. T. Valente, "Why we refactor? confessions of github contributors," in *Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering, FSE 2016*, (New York, NY, USA), pp. 858–870, ACM, 2016.
- [29] F. Calefato, L. Quaranta, and F. Lanubile, "A lot of talk and a badge: An empirical analysis of personal achievements in github," *arXiv preprint arXiv:2303.14702*, 2023.
- [30] J. Gehring, "Deterministic automatic refactoring at scale," in *2023 IEEE International Conference on Software Maintenance and Evolution (ICSME)*, pp. 541–546, IEEE, 2023.
- [31] J. White, S. Hays, Q. Fu, J. Spencer-Smith, and D. C. Schmidt, "Chatgpt prompt patterns for improving code quality, refactoring, requirements elicitation, and software design," in *Generative AI for Effective Software Development*, pp. 71–108, Springer, 2024.
- [32] A. Shirafuji, Y. Oda, J. Suzuki, M. Morishita, and Y. Watanobe, "Refactoring programs using large language models with few-shot examples," *arXiv preprint arXiv:2311.11690*, 2023.
- [33] K. DePalma, I. Miminoshvili, C. Henselder, K. Moss, and E. A. AlOmar, "Exploring chatgpt's code refactoring capabilities: An empirical study," *Expert Systems with Applications*, vol. 249, p. 123602, 2024.
- [34] Y. Gao, X. Hu, X. Yang, and X. Xia, "Context-enhanced llm-based framework for automatic test refactoring," *arXiv preprint arXiv:2409.16739*, 2024.
- [35] J. Choi, G. An, and S. Yoo, "Iterative refactoring of real-world open-source programs with large language models," in *International Symposium on Search Based Software Engineering*, pp. 49–55, Springer, 2024.
- [36] D. Wu, F. Mu, L. Shi, Z. Guo, K. Liu, W. Zhuang, Y. Zhong, and L. Zhang, "ismell: Assembling llms with expert toolsets for code smell detection and refactoring," in *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering*, pp. 1345–1357, 2024.
- [37] R. Ishizue, K. Sakamoto, H. Washizaki, and Y. Fukazawa, "Improved program repair methods using refactoring with gpt models," in *Proceedings of the 55th ACM Technical Symposium on Computer Science Education V. 1*, pp. 569–575, 2024.
- [38] D. Gautam, S. Garg, J. Jang, N. Sundaresan, and R. Z. Moghadam, "Refactorbench: Evaluating stateful reasoning in language agents through code," in *NeurIPS 2024 Workshop on Open-World Agents*, 2024.
- [39] D. Cui, Q. Wang, Y. Zhao, J. Wang, M. Wei, J. Hu, L. Wang, and Q. Li, "One-to-one or one-to-many? suggesting extract class refactoring opportunities with intra-class dependency hypergraph neural network," in *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis*, pp. 1529–1540, 2024.
- [40] D. Cui, J. Wang, Q. Wang, P. Ji, M. Qiao, Y. Zhao, J. Hu, L. Wang, and Q. Li, "Three heads are better than one: Suggesting move method refactoring opportunities with inter-class code entity dependency enhanced hybrid hypergraph neural network," in *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering*, pp. 745–757, 2024.
- [41] Y. Zhang, Y. Li, G. Meredith, K. Zheng, and X. Li, "Move method refactoring recommendation based on deep learning and llm-generated information," *Information Sciences*, vol. 697, p. 121753, 2025.
- [42] X. Gao, Y. Xiong, D. Wang, Z. Guan, Z. Shi, H. Wang, and S. Li, "Preference-guided refactored tuning for retrieval augmented code generation," in *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering*, pp. 65–77, 2024.
- [43] B. Zhang, P. Liang, Q. Feng, Y. Fu, and Z. Li, "Copilot-in-the-loop: Fixing code smells in copilot-generated python code using copilot," in *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering*, pp. 2230–2234, 2024.
- [44] Z. Zhang, Z. Xing, X. Ren, Q. Lu, and X. Xu, "Refactoring to pythonic idioms: A hybrid knowledge-driven approach leveraging large language models," *Proceedings of the ACM on Software Engineering*, vol. 1, no. FSE, pp. 1107–1128, 2024.
- [45] W. J. Conover, *Practical nonparametric statistics*, vol. 350. John Wiley & Sons, 1998.
- [46] N. Cliff, "Dominance statistics: Ordinal analyses to answer ordinal questions," *Psychological Bulletin*, vol. 114, no. 3, p. 494, 1993.
- [47] J. Romano, J. Kromrey, J. Coraggio, and J. Skowronek, "Appropriate statistics for ordinal level data," in *Proceedings of the Annual Meeting of the Florida Association of Institutional Research*, pp. 1–3, 2006.
- [48] <https://chat.openai.com/share/e9f4664e-50a5-40c1-8604-befe89a2dd36>.
- [49] <https://chat.openai.com/share/dbd163f7-aa90-4351-a9a7-a2deb906120e>.
- [50] L. Reynolds and K. McDonell, "Prompt programming for large language models: Beyond the few-shot paradigm," p. 10, 2021.
- [51] <https://chatgpt.com/share/73e56b34-fb0d-4056-bfcc-daf800b5d213>.
- [52] L. Chen, M. Zaharia, and J. Zou, "How is chatgpt's behavior changing over time?," *Harvard Data Science Review*, vol. 6, no. 2, 2024.
- [53] D. Pomian, A. Bellur, M. Dilhara, Z. Kurbatova, E. Bogomolov, T. Bryksin, and D. Dig, "Together we go further: LLMs and ide static analysis for extract method refactoring," *arXiv preprint arXiv:2401.15298*, 2024.
- [54] R. Tufano, A. Mastropaolo, F. Pepe, O. Dabić, M. Di Penta, and G. Bavota, "Unveiling chatgpt's usage in open source projects: A mining-based study," *arXiv preprint arXiv:2402.16480*, 2024.
- [55] E. Murphy-Hill, C. Parnin, and A. P. Black, "How we refactor, and how we know it," *IEEE Transactions on Software Engineering*, vol. 38, pp. 5–18, Jan 2012.