

# CoReQA: Uncovering Potentials of Language Models in Code Repository Question Answering

Jialiang Chen<sup>†§¶</sup>, Kaifa Zhao<sup>‡§¶</sup>, Jie Liu<sup>||\*</sup>, Chao Peng<sup>||</sup>, Jierui Liu<sup>||</sup>, Hang Zhu<sup>||</sup>,  
Pengfei Gao<sup>||</sup>, Ping Yang<sup>||</sup>, Shuiguang Deng<sup>†\*</sup>

<sup>†</sup> Zhejiang University, <sup>‡</sup> The Hong Kong Polytechnic University, <sup>||</sup> ByteDance

**Abstract**—Large language models that enhance software development tasks, such as code generation, code completion, and code question answering (QA), have been extensively studied in both academia and the industry. The models are integrated into popular intelligent IDEs like JetBrains and Cursor. Current benchmarks for evaluating models’ code comprehension capabilities primarily focus on code generation or completion, often neglecting QA, which is a crucial aspect of understanding code. Existing code QA benchmarks are derived from code comments with predefined patterns (e.g., CodeQA) or focus on specific domains, such as education (e.g., CS1QA). These benchmarks fail to capture the real-world complexity of software engineering and user requirements for understanding code repositories.

To address this gap, we introduce CoReQA, a benchmark for Code Repository-level question answering, constructed from GitHub issues and comments from 176 popular repositories across four programming languages. Repository-level QA requires models to retrieve the relevant content from the repository and generate answers for questions. Since questions and answers may include both natural language and code snippets, traditional evaluation metrics such as BLEU are inadequate for assessing repository-level QA performance. Thus, we provide an LLM-as-a-judge framework to evaluate QA performance from five aspects. Based on CoReQA, we evaluate the performance of three baselines, including two short-context models using generic retrieval strategies and one long-context model that utilizes the entire repository context. Evaluation results show that state-of-the-art proprietary and long-context models struggle to address repository-level questions effectively. Our analysis highlights the limitations of language models in assisting developers in understanding repositories and suggests future directions for improving repository comprehension systems through effective context retrieval methodologies.

**Index Terms**—Repository-level dataset, Code Benchmark

## I. INTRODUCTION

Large language models (LLMs) [1]–[4] demonstrate exceptional ability in understanding and processing natural language. Recently, LLMs trained on extensive code datasets (Code-LLMs) [5]–[8] show promising results in software engineering tasks. The tasks include code completion [9]–[11], code summarization [12], bug fixing [13]–[15], and et.al.

To alleviate burdens of the tasks in daily development, IDEs such as Cursor [6], Copilot [16], and JetBrains [17] integrate Code-LLMs into their features. The features include automatic

code generation within code files, bug detection and repair suggestions in the terminal, and answering user questions in chat interfaces, among others. Therefore, evaluating LLM capabilities is crucial, as it ensures that the models provide accurate, complete, and readable answers essential for software development. Rigorous evaluation helps identify the strengths and limitations of LLMs, guiding further improvements and ensuring that the models can effectively handle complex and real-world scenarios in software repositories. However, existing evaluations predominantly focus on code generation [18]–[21] or bug fixing [22], with limited attention given to code question answering.

While several recent evaluation studies introduce code QA datasets [23]–[25], the datasets primarily focus on simple yes-or-no question, and are limited method-level [24], [26] or file-level [23] context. The datasets do not reflect real-world scenarios that require understanding entire projects or clarifying inter-function dependencies, which are complex and large-scale in nature. Responses in repository-level QA require natural language explanations accompanied by relevant code examples. Existing repository-level benchmarks are limited to bug fixing [22] and code generation [18] tasks. The literature still lacks a repository-level code QA benchmark.

To address this gap, we introduce CoReQA, a benchmark designed for repository-level question answering. CoReQA is crafted to reflect the complexity and diversity of real-world inquiries, facilitating a more precise evaluation of QA systems’ capabilities within extensive code repositories. To establish CoReQA as a credible benchmark, we develop an GitHub issue-based QA generation framework. We collect real-world GitHub issues and curate a dataset of 1,563 QA pairs from 190 repositories across four programming languages. We use LLMs to reformulate issue titles and descriptions into questions, making the issues more suitable for a question-answering context. To generate reference answers, we design a prompt and require language models to understand the semantics of the question, the raw issue content, and the relevant comments. Fig. 1 gives an example from CoReQA, which is extracted from a closed GitHub issue in a popular repository and includes three positive comments. GitHub issues reflect users’ requirements for using the repository, while positive comments indicate agreement from other users. Our dataset construction framework ensures that demands represented are authentic and relevant to developers.

<sup>§</sup> Equal contribution (co-first authors). Authors are listed alphabetically by last name.

<sup>¶</sup> The job was done when the author was internship at ByteDance.

<sup>\*</sup> The corresponding authors.

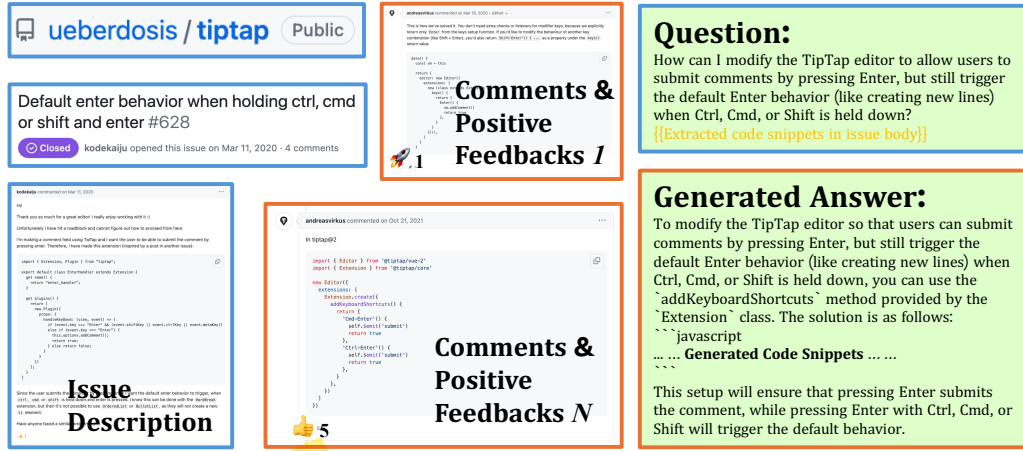


Fig. 1: A question-answer pair from a sample GitHub issue and comments.

To evaluate the effectiveness of repository-level QA, we implement a comprehensive evaluation framework using the LLM-as-a-judge [27], [28]. Our framework assesses the performance of models from two aspects: 1) absolute quality evaluation that measures accuracy, completeness, relevance, and clarity, and 2) pairwise comparison evaluation for inter-model performance assessment. We evaluate three state-of-the-art LLMs (GPT-4o [2], Gemini-1.5 [4], and DeepSeek-V2 [1]) on our CoReQA. Addressing QA pairs in CoReQA requires content from the repository. For instance, an issue may contain code snippets and pose questions, while the corresponding comments might include additional code snippets that address these questions, as illustrated in Figure 1. Thus, we evaluate the baselines under three settings: 1) query the models without any contents, 2) query the models with general retrieval strategies, i.e. BM25 [29], and 3) provide the models the contents of the whole repository. We use the third setting to evaluate long-context models on repositories where the text length fits within the model’s input token limit, such as 10 million tokens for Gemini-1.5.

Experiments demonstrate that models achieve limited success on CoReQA without additional content, with average scores of 6.37, 5.70, 7.33, and 8.09 out of 10 in accuracy, completeness, relevance, and clarity, respectively. When relevant context is provided, scores improve to 6.40, 5.74, 7.36, and 8.11 in accuracy, completeness, relevance, and clarity. The scores improved slightly with retrieved relevant context, demonstrating the benchmark’s challenge and the need for enhanced retrieval strategies to improve model performance. Additionally, even when providing the entire repository content to long-context models, performance improves only slightly, and the metrics remain suboptimal. Our results underscore the need for further research to develop effective content retrieval strategies.

Our contributions are summarized as follows:

- We provide an automated framework to construct repository-level QA pairs and a novel evaluation infrastructure to assess the performance of generative models

in addressing repository-level questions.

- We introduce CoReQA, a benchmark for repository question answering derived from GitHub repositories across various domains and programming languages. The questions are sourced from real-world user issues, with answers from corresponding discussions and comments.
- We conduct a comprehensive evaluation and comparison on state-of-the-art LLMs on CoReQA. Our evaluations include short-context LLMs with a general retrieval strategy and long-context LLMs for cross-file challenges. Evaluation results highlight the ongoing difficulty of repository-level QA tasks for LLMs and the urgent need for precise context retrieval strategies.

**Roadmap.** The remainder of the paper is organized as follows: §II provides basic background on code comprehension and retrieval argument generation systems, and formulates repository-level QA task. §III introduces the framework of CoReQA. §IV presents the experimental settings and results of baseline models. §V reviews related work. §VI discusses threats to validity, ethical considerations, limitations, and future work. §VII concludes this work.

## II. BACKGROUND

This section provides background on LLM-based code comprehension, presents prevalent strategies (e.g., retrieval-augmented generation and long-context learning) for solving repository-level code understanding tasks, and formally define the repository-level question answering task.

### A. LLM-based code comprehension

In the software development, code comprehension proves indispensable [30]–[32] due to the iterative dependencies among code elements that collectively enable intricate functionalities. For instance, to automatically generate a code segment [19], [33], developers need to understand the imported modules and discern between identically named functions belonging to different classes. Similarly, crafting an accurate pull request [22] for a repository necessitates a clear grasp

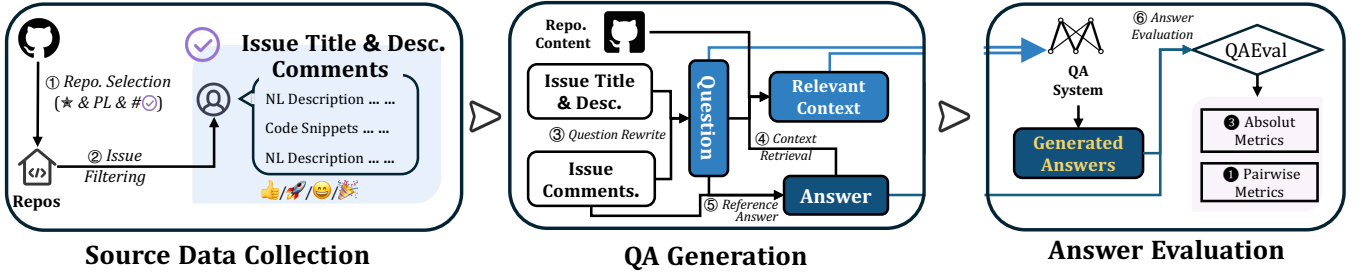


Fig. 2: The overall pipeline of CoReQA.

of requirements, alongside the aggregation of pertinent data, logs, and error messages.

Despite the advancements in large language models [1]–[4], effectively tackling cross-file code comprehension remains a substantial hurdle. Interpreting the intent behind cross-file requirements is often inconsistent due to inherent ambiguity [34], [35]. Additionally, retrieving relevant contents scattered across different files is challenging [28], [33]. Constructing the appropriate context to identify the core issue and develop a resolution is also a complex task. The challenges underscore the urgent need for advancements in models and methods that can effectively navigate and understand the intricate interdependencies within **code repositories**.

### B. RAG and long-context learning

Retrieval-augmented generation (RAG) enhances natural language processing (NLP) by integrating retrieval-based techniques with generative models. RAG addresses the limitations of LLMs that often lack access to real-time and comprehensive external information. RAG framework combines a retriever component that fetches relevant content from a large corpus with a generator component that synthesizes this information into coherent responses [36]. The retriever identifies pertinent content based on the input query, while the generator uses the content to craft a detailed and accurate response [37], [38]. The combination of retriever and generator enhances the utility of LLMs in applications requiring up-to-date or specific knowledge. In the domain of code repositories, the RAG framework retrieves relevant code snippets from various files and generates comprehensive explanations or solutions [25], [33]. The RAG workflow involves processing a user query, retrieving pertinent documents, and generating a synthesized answer, ensuring that the response is both accurate and informative [24]. The RAG workflow enhances the accuracy of responses by integrating real-time and relevant information [21], [22], [24]. The approach leverages dynamic and external data that static models may miss.

Long-context learning is an emerging solution [39]–[41] for repository-level code comprehension, a methodology renowned for its capacity to analyze the entirety of a repository’s content. Nonetheless, this approach encounters limitations when repositories become excessively large, surpassing the feasible boundaries of what long-context models can efficiently process.

### C. Repository-level question answering task

In this work, we focus on constructing a benchmark and setting up baselines for repository QA tasks. Focusing on a target code repository  $\mathcal{R}$  and confronted with a repository-related question  $q$ , which comprises natural language, function or class signatures, or even code snippets, the task necessitates the production of a free-form textual answer  $a$ . The textual answer  $a$  is versatile, encompassing the possibility of being a sentence, a phrase, or even a code snippet, contingent upon the context and requirements of the query. Indeed, distilling the answer  $a$  from a single file would be overly simplistic and inadequate given that  $q$  aims to probe or leverage insights about the entire repository  $\mathcal{R}$ . Consequently, formulating a satisfactory response  $a$  necessitates a deep comprehension of the repository’s architecture alongside the retrieval and synthesis of pertinent information scattered throughout  $\mathcal{R}$ . Therefore, the repository-level question-answering task is a difficult challenge, which is exacerbated by the inherent scalability issues of large code bases and the complexity of long-term context understanding. It is of immense significance and urgency to establish and disseminate a rigorous, representative benchmark for this domain.

## III. CoReQA

In this section, we introduce the construction process of the CoReQA as illustrated in Fig.2. The construction of the CoReQA encompasses raw data collection, question-answer (QA) pair generation, and QA evaluation.

### A. CoReQA construction

1) *Raw Data Collection*: We design following pipelines to make CoReQA authentic, reliable, and comprehensive.

**Repository Selection.** We start by selecting the top 500 GitHub repositories, ranked by the number of stars and having valid licenses [42] across four programming languages: Python, Java, Go, and TypeScript. Next, we filter out repositories where the number of closed issues and pull requests is either lower than 1,000 or higher than 50,000. The selection ensures the repositories are not only popular but also reflective of active development communities in these prominent programming languages. We set a maximum threshold of 50k to reduce the analysis burden because every time we analyze the issue, we need to request GitHub API once. We obtain 890 repositories.

```

reference_answer_generation_prompt = """\
## GitHub Issue Answer Generator 🚀

You are an AI assistant designed to help summarize GitHub issue
comments and generate helpful answers for repo-related questions.
Given a GitHub issue title, body, comments, and a concise user
query, follow these steps to generate a clear and concise answer:

1. Read the issue title, body, and comments to understand the
core problem or questions.
2. Summarize the key points from the comments and extract the
most relevant information to generate a helpful answer to the
user query.
3. Provide a clear, actionable solution or recommendation that
directly addresses the issue.
4. If the issue contains code snippets, you can refer to them in
the answer, but avoid including long code segments.
5. If there are no feasible solutions in the comments, Simply
response "I am sorry, I could not find a solution for this
issue."
5. Directly output the answer without any additional information
or context. No greetings, No thanks, No closing remarks.
6. Do not include any personal information or specific user
details in the answer, do not mention the source of the
information, specially, don't mention this solution from issue.

## Real Issue:
### Repository Name
{repo_name}
### Issue Title
{issue_title}
### Issue Body
{issue_body}
### Relevant Information from Issue Comments
{issue_comments}
### Question to be answered
{summarized_question}

### Summarized Issue Answer:
"""

```

Fig. 3: Reference answer generation prompt.

To ensure that CoReQA is capable of evaluating long-context models, token lengths for each repository are estimated using OpenAI’s tokenizer calculation method [43]: the token length is approximated by dividing the number of characters by four. We then select 50 repositories with token lengths greater than 200K and 5 with less than 200K [3], [4] for each programming language. To this end, we obtain 218 candidate repositories with 1,623,624 issues.

**Issue Collection and Filtering.** For candidate repositories, all closed issues, along with their accompanying comments (Issue-with-Comments, IwC), are crawled. We focus on identifying issues suitable for conversion into question-answering pairs. Firstly, we filter out issues that do not contain the tags “feat,” “bug,” “fix,” or “error.” Our human inspection finds that the aforementioned tags typically denote requests for code fixes or pull requests, which do not usually have direct solutions in the issue comments. For this study, we evaluate text-only models. We exclude issues and comments containing images, as these images are often screenshots of code or terminal logs that text-only models cannot process. We concentrate on issues related to code descriptions and select those with descriptions that include code snippets. We select issues with at least three positive comments to ensure the comments help address the problem and generate high-quality

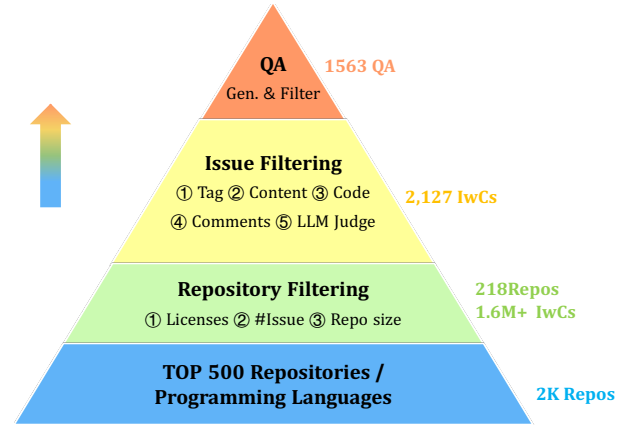


Fig. 4: Filtering and selection process for CoReQA.

question-answering pairs. Positive comments are those marked with “+1,” “laugh,” “hooray,” “heart,” and “rocket” from other users. Additionally, some comments may direct users to refer to other closed issues or commit to resolving the current issue. We exclude issues if the positive comments provide such links. To this end, we obtain candidate 8,977 issues from 213 repositories. To ensure QA pairs are general and cover a broader range of repositories, we select a maximum of 300 issues with comments from each candidate repository. Furthermore, to make these issues suitable for generating QA pairs, we leverage the semantic analysis capabilities of LLMs and design a specific prompt to enable the LLM to pre-filter the issues effectively. To this end, we obtained candidate 2,127 issues from 190 repositories.

2) *QA pair generation*: Question-answer pairs in CoReQA are meticulously derived from the GitHub repository’s issue titles, descriptions, and associated comments, ensuring relevance to real-world development scenarios.

**Question Generation.** Questions in CoReQA are derived from real-world GitHub issues. We use LLMs to rewrite questions with selected issues with positive comments. We apply prompt engineering to design instructions that guide LLMs to analyze issues and rewrite the issues as repository-related questions. Specifically, we use *Self-Consistency* [44] to format the LLM’s output, ensuring that the generated questions are easily extracted and analyzed. Additionally, we apply *Chain-of-Thought* (CoT) to enable the LLM to thoroughly analyze the issues and produce more relevant questions. Recognizing that the issue contents might be novel to the generative model, our prompt ingeniously integrates both the issue title and description. Since issues often contain code snippets that are crucial for addressing the question, we also attach the original code snippets extracted from the issue body to the final question. Furthermore, to uphold the relevance and value of generated questions, we enforce a rigorous in-house human inspection process [45], thereby ensuring a high quality.

**Reference Answer Construction.** We devise a prompt to generate reference answer integrating the question, issue’s tag, issue title, issue description, and corresponding comments.



TABLE I: Statistics of CoReQA.

Language	# Repositories	# QA Pairs
Python	46	439
Java	33	209
Go	48	371
TypeScript	49	544
Total	176	1,563

Fig. 3 illustrates the prompt for reference answer construction. The yellow strings enclosed in curly braces indicate content that will be replaced depending on the specific repository and issue information. First, we assign a role to the LLM to generate answers. Subsequently, CoT prompt engineering is employed to instruct language models in step-by-step understanding and analysis of issues and their corresponding comments. We also specify the requirements for the LLM to format its output properly. To this end, we obtain 1,563 candidate question-answering pairs, specifically, 44 of these questions support long-context model evaluation. Fig.4 sketches the data selection process for CoReQA.

**Related Content Retrieval.** CoReQA provides reference content extracted from the repository to assist in answering questions. QA pairs in CoReQA are meticulously derived from real-world GitHub Issues with Comments (IwCs), which inherently contain critical context, including relevant code snippets and detailed descriptions. When users pose questions within an AI IDE or a QA platform, the models only have access to the repository’s source code. To assess QA system’s capacity to generate appropriate answers, CoReQA provides related contents from the repository that are contextually relevant to each question. To achieve this, we use LangChain [46] to split files in the repository into text chunks. BM25 [22], [29], [47] is then employed to retrieve relevant chunks from the entire repository based on two sources: the code snippets in the issue body and the generated questions. The top 5 chunks from each source are selected, resulting in 10 reference content for each QA pair. It should be noticed that we apply a generic strategy to retrieve relevant content, the content is provided as reference materials for LLM to address questions. Furthermore, in order to evaluate Repo-as-text, we organize the entire contents of the repository using Markdown format as final context. The statistics for CoReQA is given in Table I.

### B. Components of CoReQA

Each QA pair in the CoReQA consists of the following information:

**Issue information.** We provide the issue’s URL, issue title, issue descriptions, host repository name, and all issue comments with feedback tags.

**Question.** The questions are in a mixed format of natural language description, and code snippets in Markdown format.

**Answers.** The answer may only include natural language response, code snippets, or the mixed format of natural language and code snippets.

**Reference context.** Ten retrieval content obtained using BM25.

### C. CoReQA evaluator

CoReQA evaluator utilizes the LLM-as-a-judge approach [27], [28], and focuses on two key aspects: absolute quality evaluation and pairwise comparison evaluation. We design prompts<sup>2</sup> to evaluate the absolute quality from four dimensions. To mitigate the impact of scoring temperature, each absolute value scoring is performed twice, with the average taken as the final result. To verify the stability of the scoring process, an additional experiment was conducted: 200 question-answer pairs (QAs) were randomly selected from the sample, and each pair was scored five times. The mean score and variance were then calculated for each QA pair. Given the variability in LLM-generated answers, the CoReQA Evaluator reports the average and standard deviation for these metrics over five trials to ensure a robust evaluation. The detailed scoring criteria are as follows:

**Accuracy (Acc.)** assesses the factual correctness of the generated answer in comparison with the reference answer. We employ the chain-of-thought (CoT) [48] strategy to compare the generated answer with the reference answer. The judge verifies the factual correctness of the generated answer and checks the accuracy of quoted sources. The judge rates the generated answer in comparison to the reference answer on a scale of 1 to 10, categorized into five levels. In detail, a score of 1-2 means the generated answer is mostly incorrect; 3-4 denotes that the generated answer contains significant factual errors; 5-6 indicates that the generated answer has some factual errors but is primarily accurate; 7-8 means the generated answer has minor inaccuracies but is overall correct; and 9-10 denotes that the generated answer is factually correct and has no errors.

**Completeness (Cmpt.)** evaluates whether the model’s answer covers all aspects of the question. We employ the CoT prompt strategy to instruct the judge to understand the key components of the reference answer and identify any critical points absent in the generated answer. We require the judge to rate the generated answer on a scale of 1 to 10 across five levels.

**Relevance (Rel.)** assesses whether the generated answer addresses the core concern of the question. The judge is asked to understand the core concern of the question, determine whether the generated answer directly addresses the question, and make sure the generated answer stays on-topic. The judge rated the generated answer regarding the question on a scale of 1 to 10 across five levels.

**Clarity** determines whether the generated answer is easily understandable. The LLM judge evaluates the logical clarity and simplicity of the generated answer, ensuring it is easy to comprehend. The LLM judge rates the generated answer in comparison to the reference answer on a scale of 1 to 10, categorized into five levels.

For each evaluation metric, the performance of the models will be rated on a scale from 1 to 10, with a higher score indicating superior performance.

Pairwise comparison evaluation (**PCE**) compares the generated answers from two models against a reference answer to

TABLE II: Evaluation model details and configuration settings.

Model	Context Length	Temperature
GPT-4o	8,192	0.2
DeepSeek-V2	128k	0.2
Gemini-1.5	10M	0.2

determine which model provides a superior response. To mitigate the impact of answer order on the results, the CoReQA Evaluator conducts PCE twice: once with one model’s answer presented first, and once with the other model’s answer presented first. Besides, we use the Elo score as the dual evaluation metric [49], [50]. The Elo score enables the ranking of models based on their performance in comparative evaluations.

The dual evaluation framework, i.e., absolute quality evaluation and pairwise comparison evaluation, ensures a comprehensive assessment of AI-generated answers by balancing both absolute and relative scoring methods. The absolute scoring evaluates each answer independently based on predefined quality criteria, while the relative scoring directly compares the answers against each other.

#### IV. EXPERIMENTS

##### A. Experimental setup

In this section, we outline the experimental setup to evaluate the effectiveness of language models. The assessment evaluations encompass both short-context models and long-context models, employing rigorous methodologies to ensure comprehensive and reliable evaluations. The performance of short-context models is investigated under two settings: with and without the integration of retrieval contents. The setting allows for the assessment of the impact of supplementary context from the repository on the models’ question-answering capabilities. For long-context models, the entire repository context is concatenated, and questions are contextualized within this extended narrative, testing the models’ abilities to parse and utilize extensive repository information for accurate responses. To mitigate variability in answer generation and reduce token consumption, each experiment is executed five times. Statistical analyses of the results are then conducted to measure deviations, ensuring the reliability and reproducibility of our findings.

1) *Short-context models*: We consider the following widely recognized LLMs for evaluation, including GPT-4o [2]<sup>1</sup>, Gemini-1.5 [4] and DeepSeek-V2 [1].

2) *Long-context models*: Long-context models are specifically designed to accommodate and comprehend extended sequences of text, thereby transcending the context window limitations inherent in foundational models. This enhancement empowers language models to assimilate a wealth of information concurrently during inference, a critical capability for tasks that demand a broad sweep of contextual understanding—characteristic of the challenges encountered in the CoReQA scenario. We consider Gemini 1.5 [4], a

<sup>1</sup>The GPT-4o model used is procured from Azure.

TABLE III: Overall results.

Context	Models	Acc.	Cmpt.	Rel.	Clarity	PCE
None	GPT-4o	<b>6.85</b>	<b>6.27</b>	<b>7.81</b>	<b>8.40</b>	<b>61.9</b>
	DeepSeek-V2	6.41	5.77	7.39	8.19	50.0
	Gemini-1.5	5.86	5.05	6.79	7.69	37.6
BM25	GPT-4o	<b>6.91</b>	<b>6.34</b>	<b>7.86</b>	<b>8.42</b>	<b>57.9</b>
	DeepSeek-V2	6.39	5.79	7.39	8.18	50.0
	Gemini-1.5	5.90	5.08	6.83	7.73	33.2

popular large language model (LLM) for evaluation, which can handling 10M token context.

Table II provides the basic information of the models under evaluation, including the context length and temperature settings during the evaluation. Across all models under evaluation, we adopt a consistent setting for the temperature parameter at 0.2. This configuration encourages the generation of more deterministic and focused responses, thereby promoting stability and reliability in the inference outcomes. Furthermore, we adhere to a single-round instruction approach to streamline the interaction process and ensure that the evaluation is standardized, leading to robust and comparable results across different models.

3) *Research Question*: We design experiments to address the following research question:

**RQ1: Main Results.** What is the efficacy of models in addressing questions within the CoReQA?

**RQ2: CoReQA Investigation.** How do LLMs perform across QA pairs with different properties?

**RQ3: QA Evaluator Validity.** Does the LLM-as-a-judge based evaluator effective in measuring question answering performance?

B. *RQ1: What is the efficacy of models in addressing questions within the CoReQA?*

Table V presents the overall results of our evaluation, comparing different models across two settings on all QA pairs in CoReQA: 1) addressing questions without any reference contents, and 2) answering questions with BM25 retrieval contents. For pairwise comparison evaluation, we selected DeepSeek-V2 as the short-context model for comparison [49] due to its competitive performance in code-related tasks and its open-source availability [1]. Given the high cost of model inference tokens, the expense of comparing models pairwise for metric calculation is unacceptable. For absolute quality evaluation, Table V demonstrates that GPT-4o outperforms all other models under both no-context and BM25 retrieval context settings, while Gemini-1.5 shows relatively poor performance. It can be observed that all models achieve scores of 5-6 for completeness (Cmpt.), indicating that the models cannot fully address all aspects of the problems, regardless of reference context. Except for Gemini-1.5, all models achieve desirable clarity, scoring 7-8, which means the answers are mostly clear and easy to understand. For accuracy (Acc.) and relevance (Rel.), three models score between 6 and 8, indicating that the generated answers have few errors and are

**TABLE IV: Model performance on QA pairs across different programming languages.**

		Go					Java				
	Models	Acc.	Cmpt.	Rel.	Clarity	PCE	Acc.	Cmpt.	Rel.	Clarity	PCE
None	GPT-4o	6.69	6.13	7.64	8.38	63.2	6.65	6.21	7.68	8.29	55.4
	DeepSeek-V2	6.18	5.63	7.22	8.12	50.0	6.34	5.80	7.36	8.20	50.0
	Gemini-1.5	5.84	5.06	6.72	7.72	38.2	5.83	5.08	6.81	7.64	38.0
	<i>Average</i>	6.24	5.61	7.19	8.07	/	6.27	5.70	7.28	8.04	/
BM25	GPT-4o	6.73	6.25	7.70	8.39	58.3	6.71	6.21	7.67	8.25	54.2
	DeepSeek-V2	6.17	5.63	7.16	8.11	50.0	6.32	5.83	7.32	8.11	50.0
	Gemini-1.5	5.81	5.06	6.70	7.73	36.1	5.94	5.19	6.95	7.71	32.2
	<i>Average</i>	6.24	5.65	7.19	8.08	/	6.32	5.74	7.31	8.02	/
		Python					TypeScript				
None	GPT-4o	7.00	6.35	7.93	8.43	65.7	6.92	6.32	7.89	8.43	60.8
	DeepSeek-V2	6.48	5.79	7.48	8.22	50.0	6.49	5.84	7.44	8.21	50.0
	Gemini-1.5	5.79	4.99	6.78	7.65	36.0	5.93	5.08	6.84	7.72	38.8
	<i>Average</i>	<u>6.42</u>	<b>5.71</b>	<b>7.40</b>	<u>8.10</u>	/	<b>6.45</b>	<b>5.75</b>	<u>7.39</u>	<b>8.12</b>	/
BM25	GPT-4o	7.03	6.38	7.95	8.45	58.4	7.01	6.42	7.95	8.46	58.5
	DeepSeek-V2	6.53	5.86	7.51	8.19	50.0	6.51	5.84	7.47	8.24	50.0
	Gemini-1.5	5.84	4.97	6.79	7.71	30.4	5.98	5.15	6.89	7.77	33.9
	<i>Average</i>	<u>6.47</u>	<u>5.74</u>	<u>7.42</u>	<u>8.12</u>	/	<b>6.50</b>	<b>5.80</b>	<b>7.44</b>	<b>8.16</b>	/

mostly relevant to the questions. Pairwise comparison results (PCE) also show that GPT-4o generates better results than DeepSeek-V2 (61.9 PCE vs. 50 PCE), and Gemini-1.5 shows significantly worse results. It can be observed that even though DeepSeek-V2 and Gemini-1.5 can achieve comparatively good performance on clarity, those models cannot infer complete answers for questions and miss a few critical points. The improvement brought by BM25 is also limited, which hints at adjustments to the ability to retrieve high-quality relevant information to improve the model QA capability.

Comparing the results of no-context and BM25 retrieved-context settings, all models show only limited improvements. This limited enhancement may be attributed to BM25 retrieval strategy, which might not provide sufficiently helpful content for effectively addressing the questions. Issues in the repository is often challenging, as it requires a comprehensive understanding of the overall project structure and precise control over project details to accurately locate the most relevant context. By using LangChain to split the repository into chunks and BM25 to retrieve the most similar content based on code snippets in the questions, the semantic relevance of the retrieved content may be fragmented, offering limited useful information. Additionally, since BM25 relies on term frequency to find similar content, it may struggle to capture the ideal context, as terms in code often exhibit minimal variation. Therefore, relying solely on word similarity is insufficient for this task. A more nuanced approach that considers the semantic structure and the specific context of the code is necessary to effectively address the issues.

Table V presents the results of Gemini-1.5 on 44 QA pairs from repositories where the content length is within Gemini-1.5’s capacity limitations. The results demonstrate that as the length of reference content increases, the performance of Gemini-1.5 improves correspondingly. This observation

**TABLE V: Evaluation of long-context model.**

Context	Acc.	Cmpt.	Rel.	Clarity
None	5.81	4.86	6.74	7.49
BM25	5.86	5.16	6.93	7.55
Repo	5.94	5.45	6.86	7.59

suggests that Gemini-1.5 benefits from longer contexts, enhancing model performance. Additionally, because the long-context approach outperforms the BM25 retrieval in terms of effectiveness, it indirectly indicates that BM25 may not effectively retrieve the most relevant information for QA pairs.

*Answer to RQ1:* The repository question-answering task remains a challenge for large language models. While retrieving relevant information can enhance model performance in QA scenarios, the quality of the retrieved information is crucial. The performance of long-context models improves with more extensive reference information, but it is ultimately constrained by the model’s inherent capabilities.

*C. RQ2: How do LLMs perform across QA pairs with different properties?*

a) *QA pairs across different programming languages:* Table IV illustrates LLMs performance on QA pairs in CoReQA across different programming languages and gives the highest score in bold font and second highest score with underline. Table IV shows that models achieve the best overall performance on TypeScript-related QA pairs, followed by those related to Python. The phenomenon could lie in the property of program languages that TypeScript provides more consistent syntax and structure than other languages, while Python works as an interpreted language and provides

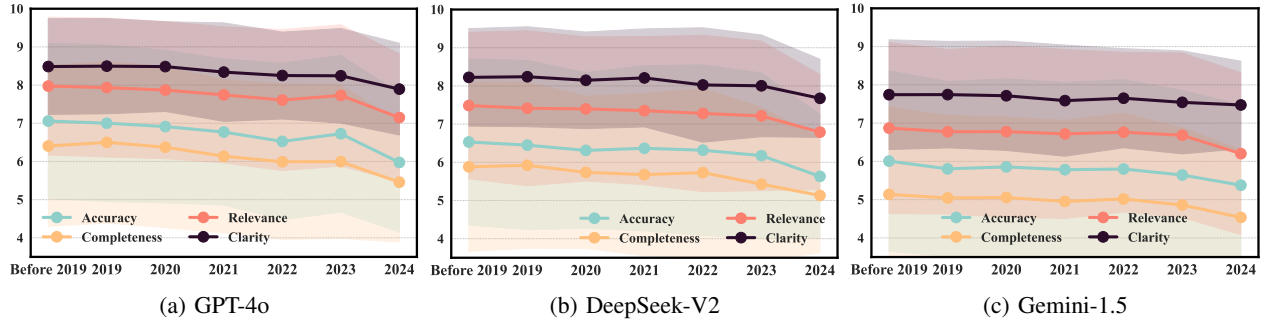


Fig. 5: Performance of models on QA pairs over various time periods.

precise grammar. GPT-4o demonstrates higher PCE values than other models under all settings in all programming languages, while Gemini-1.5 performs worse. Table IV also shows that providing relevant content to address the questions narrows the performance gap between GPT-4o and DeepSeek-V2. This phenomenon may be attributed to GPT-4o’s superior ability to leverage data already present in model’s training set, whereas DeepSeek-V2 is more effective in utilizing the reference information provided in the prompt.

b) *QA pairs across different timelines*: Fig. 5 illustrates the performance of different models on QA pairs across timelines based on the creation dates of the source issues associated with the QA pairs. The variance for each period is notably high, with the shaded area in each color representing an average variance of approximately 1.5. This indicates a wide range of difficulty levels among the QA pairs in CoReQA. Additionally, the model’s performance on QA tasks declines over time. The decline correlates with the recency of the model’s training data; specifically, performance drops when issues are less recent relative to the training data. The observation indicates that the models perform better on issues that are potentially in the training data of models. Furthermore, it suggests that user questions become more complex over time, necessitating a deeper understanding of the repository’s functionality for effective problem resolution.

c) *QA pairs with vary token lengths*: Fig. 6 shows the performance of GPT-4o across different token length ranges for questions: 37-118, 119-184, 185-286, 287-538, and 539-3540. We divide the question lengths into quintiles to determine the ranges. The evaluation is conducted without any reference contents. Models perform best on questions with token lengths in the 287-538 range, achieving the highest scores in accuracy, completeness, and relevance. This may be because questions of this length provide sufficient information, such as relevant code snippets, without overwhelming the model and causing it to lose focus on the question. If a question’s length is too short, such as within the range of 37-118 tokens, it may lack sufficient detail, which can reduce the model’s ability to address the issue effectively. Conversely, if a question is too long, exceeding 539 tokens, the model may struggle to retain key information, potentially impairing its performance. Notably, query length significantly impacts the Completeness

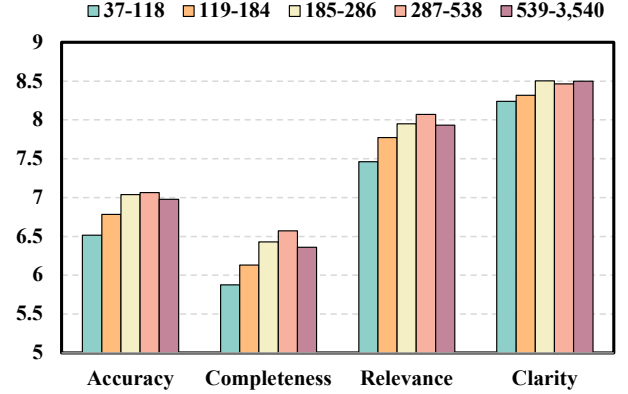


Fig. 6: Performance of GPT-4o over various question token length.

and Relevance metrics, with score differences ranging from 0.7 to 0.8 points between the highest and lowest values, indicating that the model’s ability to deliver comprehensive and relevant responses is susceptible to the length of the input prompt. In contrast, the Accuracy and Clarity metrics show less variability, with differences of approximately 0.5 points, suggesting that the model’s accuracy and clarity are relatively unaffected by changes in prompt length.

*Answer to RQ2:* Language models are sensitive to properties of QA pairs. The models perform better with well-structured programming languages and are more effective with QA pairs where the data source predates the model’s training data. The length of QA pairs impacts performance; excessively long questions and those with insufficient information impair models’ effectiveness.

D. *RQ3: Does the LLM-as-a-judge based evaluator effective in measuring question answering performance?*

This research question evaluates the effect of the CoReQA evaluator in our scenario. We first assess the robustness of the absolute quality evaluator. We randomly select 200 question-answering pairs in CoReQA and the corresponding generated answers by each model and then evaluate the absolute quality evaluator five times. Table VI gives the average and standard deviation. The results demonstrate that the absolute quality



TABLE VI: The validity of absolute quality evaluation based on the average and standard deviation.

Strategies	Models	Acc.	Cmpt.	Rel.	Clarity	BLEU
None	GPT-4o	7.14	6.56	8.03	8.54	7.7
		( $\pm 0.59$ )	( $\pm 0.67$ )	( $\pm 0.59$ )	( $\pm 0.49$ )	
	DeepSeek-V2	6.68	6.14	7.66	8.35	8.2
		( $\pm 0.49$ )	( $\pm 0.57$ )	( $\pm 0.54$ )	( $\pm 0.49$ )	
	Gemini-1.5	6.03	5.25	6.94	7.76	4.4
		( $\pm 0.61$ )	( $\pm 0.62$ )	( $\pm 0.64$ )	( $\pm 0.55$ )	
BM25	GPT-4o	7	6.32	7.92	8.48	6.3
		( $\pm 0.58$ )	( $\pm 0.62$ )	( $\pm 0.58$ )	( $\pm 0.53$ )	
	DeepSeek-V2	6.43	5.73	7.42	8.14	5.7
		( $\pm 0.51$ )	( $\pm 0.53$ )	( $\pm 0.58$ )	( $\pm 0.52$ )	
	Gemini-1.5	6.01	5.18	6.89	7.76	5.2
		( $\pm 0.55$ )	( $\pm 0.55$ )	( $\pm 0.62$ )	( $\pm 0.54$ )	

evaluator achieves around 0.55 standard deviation five times running. The 5.5% (0.55/10) deviation is acceptable for such inference scoring tasks [27]. It does not significantly affect the scoring categories because our absolute quality evaluator requires the judge to score generated answers on a scale of 1 to 10 and categorize them into five levels with two 2-point scales each. Additionally, we provide BLEU scores for comparison. Although BLEU can differentiate the performance of different models, it does not provide distinctions based on semantic accuracy. Furthermore, in scenarios lacking reference content, DeepSeek-V2’s BLEU score is even higher than GPT-4o’s, indicating that DeepSeek-V2 tends to produce results that are more token-similar rather than necessarily more accurate, complete, relevant, or clear.

For pairwise evaluation comparison, we compare answers generated by two models twice: first with model A’s answer in front, and then with model B’s answer in front. Fig. 7 illustrates comparison results, with DeepSeek-V2 as the baseline model, consistent with previous settings. Fig. 7 shows that in 5 out of 10 cases, the judge tends to assess the answer presented first as the better one, indicating a bias towards the first position. Based on previous analysis, GPT-4o outperforms DeepSeek-V2, and DeepSeek-V2 outperforms

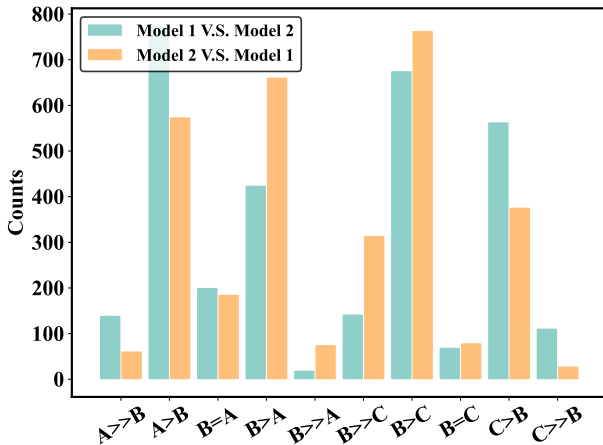


Fig. 7: Pairwise evaluation comparison of evaluation results across models. A denotes GPT-4o, B denotes DeepSeek-V2, and C denotes Gemini-1.5.

TABLE VII: Pairwise evaluator comparison

		DeepSeek-V2							DeepSeek-V2				
		≤	<	=	>	≥			≤	<	=	>	≥
GPT-4o	≤	37	73	1	22	6	Gemini-1.5	≤	22	65	1	15	8
	<	18	363	75	304	14		<	5	227	37	258	36
	=	0	42	85	72	1		=	1	13	24	30	1
	>	6	96	24	254	44		>	0	69	16	425	165
	≥	0	0	0	9	10		≥	0	2	1	35	104

Gemini-1.5. This positional bias does not affect overall judgments of the model’s capability. For instance, when comparing with DeepSeek-V2, regardless of which model’s answer is presented first, the better model is correctly identified, as shown in the corresponding bar chart (i.e.,  $A \gg B$ ,  $B \gg C$ ).

We also analyze the consistency of the pairwise comparison results from the two evaluations. Table VII shows the results of pairwise comparisons. For example, in the left table, the number in the 1st row and 3rd column indicates that when GPT-4o’s answer was presented first, the judge rates it much better than DeepSeek-V2 ( $GPT-4o \gg DeepSeek-V2$ ). When GPT-4o’s answer was presented second, the judge rates them as equal ( $GPT-4o = DeepSeek-V2$ ). Data with a muted mauve background indicate that the results of the two comparisons are consistent, regardless of which model’s answer is presented first. Data with a light pinkish background indicate a small discrepancy between the two comparisons, such as one comparison judging model A’s result as superior to model B’s, while the other comparison judges model A’s result as equivalent to model B’s. Additionally, the models rarely make significant errors when distinguishing between markedly different answers. For instance, when comparing the results of GPT-4o and DeepSeek-V2, only 6 results showed one comparison judging GPT-4o as far superior to DeepSeek-V2, while another comparison judged GPT-4o as far inferior to DeepSeek-V2, with no similar misjudgments observed. A similar conclusion can be drawn when comparing the results of Gemini-1.5 and DeepSeek-V2.

*Answer to RQ3:* The CoReQA Evaluator’s absolute quality assessments are stable and reliable, with a margin of error of 5.5%. Pairwise evaluations tend to favor answers presented first. The bias does not significantly impact cases with large performance differences between models.

## V. RELATED WORK

CodeQueries [23] is a benchmark designed to assess the capability of language models in understanding code semantics through extractive question-answering. CodeQueries includes 52 queries that necessitate single-hop and multi-hop reasoning over Python code. Each query is annotated with answer spans and supporting facts. CodeQueries presents a significant challenge for models, making it a valuable tool for advancing research in code comprehension and program analysis. CodeQA [24] constructs question-answering pairs specifically for methods within the code base. CodeQA employs template-

**TABLE VIII: Comparison of related benchmarks.**

Benchmark	Size	Task	Gran.	Language
CoderEval	460	CG	Cross file	Python, Java.
CodeAgent	101	CG	Cross file	Python.
SWE-Bench	2,294	CG	Cross file	Python.
HumanEval	164	CG	Single file	Python.
CrossCodeEval	9,928	CC	Cross file	Python, C#, Java, TypeScript.
CodeQueries	52	QA	Cross file	Python.
CodeQA	190k+	QA	Single file	Java, Python.
CS1QA	9,237	QA	Single file	Python.
CoReQA	1,563	QA	Cross file	Python, Java, TypeScript, Go.

QA: Question Answering; CG: Code Generation;  
CC: Code Completion; NL: Natural Language; Gran.:Granularity

based methods to generate question-answer pairs for Python and Java. CS1QA [25] is a benchmark designed for code-based question answering in educational contexts. It comprises 9,237 annotated question-answer pairs sourced from introductory Python programming courses. CS1QA challenges models with tasks such as question type classification, code line selection, and answer retrieval, underscoring the complexity of integrating natural language understanding with code comprehension. Initial baseline results indicate substantial room for improvement, particularly in tasks requiring detailed code analysis.

SWE-bench [22] is a benchmark designed for evaluating large language models in practical code generation and bug-fixing tasks. The benchmark utilizes publicly available pull requests from popular Python repositories. While SWE-bench offers a comprehensive evaluation framework, it is constrained by the reliance on open-source data and the current models’ context-handling limitations. CoderEval [18] is designed to evaluate code generation models in real-world settings, including both standalone and non-standalone functions from open-source projects. CoderEval addresses the gap in evaluating context-dependent code generation, providing a more comprehensive and pragmatic assessment. HumanEval [20] is constructed for evaluating code generation tasks and comprises 164 programming problems, each with a function signature, docstring, function body, and unit tests for automatic evaluation. Other related benchmarks [21], [51] are designed to address code generation tasks. While some of these benchmarks consider cross-file scenarios to better mimic real-world conditions, most are limited to the Python programming environment. We give the statistical comparison of code comprehension related benchmarks in Table VIII.

## VI. DISCUSSION

### A. Threats to validity

One threat comes from the randomness inherent in the inference process of LLMs. We set LLM’s temperature for generating diverse reference answers to 0.8. As a result, even with the same prompt and model, the generated answers may vary. Although we have set the temperature for evaluation models to 0.2 when using LLM-as-a-judge [27] to preserve the robustness of generated answers, the generated scores and

answers can still exhibit some randomness. Despite we conduct evaluations five times to mitigate this effect, variability remains. Additionally, while all prompts used for benchmark construction and evaluation have been refined through prompt engineering [52], there is no guarantee that the prompts are optimal for each LLM. With ongoing advancements in prompt engineering, better prompts may be developed that can generate better answers for different models.

### B. Ethics statement

CoReQA is constructed entirely from public code repositories with permitted license [53]. During the collection process, we do not include information about GitHub users and only collect issue title, issue descriptions and comment content. Besides, CoReQA also rewrites and reconstructs those collected information into appropriate questions answering pairs with human inspections. All annotators in CoReQA are authors of this work. To ensure the quality and consistency of the annotations, we conducted manual sampling and verification on a subset of the annotated questions and answers. CoReQA’s filtering criteria for GitHub repositories are based on popularity, measured by the number of stars, and purely random selection, ensuring that the process does not implicitly or explicitly rely on any discriminatory or biased heuristics for repository selection.

### C. Limitation and future work

CoReQA task is limited to four programming language due to constraints in human resources and token availability. We aim to extend CoReQA to cover more programming languages in the future. Since BM25 is not an effective retrieval approach to CoReQA, we will consider advanced retrieval strategies, such as applying static analysis methods as code splitters, to preserve the semantics of code snippets. In addition to these improvements, we need to consider multi-turn dialogues in future research. This study only addresses single-turn dialogues, but understanding code-related questions often requires handling more complex, multi-turn interactions. Lastly, while this work evaluates models using the LLM-as-a-judge strategy, the evaluation method relies on the inherent capabilities of the LLM. In future work, we will investigate additional linguistic and programming methods to evaluate code-related QA metrics.

## VII. CONCLUSION

In this paper, we introduce a novel benchmark named CoReQA for evaluating a model’s effectiveness in understanding code-related questions at the repository level. CoReQA is sourced from real-world code repository issues and related comments. We provide a construction pipeline to automatically expand the benchmark to different programming languages and augment its scale. Additionally, we design a comprehensive evaluation tool to assess the performance of QA system. The evaluation framework includes both absolute quality evaluation for measuring metrics from four aspects and pairwise comparison evaluation for inter-model performance

assessment. Experimental results demonstrate that answering repository-level questions remains a challenge for LLMs, even when provided with the entire repository content. We hope that our benchmark and other contributions will aid in the development of better code-related models and assist developers in building more effective platforms for understanding code repositories in the future.

## REFERENCES

- [1] D. AI, “Deepseek: Advanced ai capabilities,” <https://www.deepseek.com/en>, 2024, accessed: 2024-06-18.
- [2] OpenAI, “Chatgpt: A large language model,” <https://www.openai.com/chatgpt>, 2024, accessed: 2024-06-18.
- [3] Anthropic, “Claude: An ai assistant,” <https://claude.ai/>, 2024, accessed: 2024-06-18.
- [4] G. DeepMind, “Gemini: A multimodal ai,” <https://gemini.google.com/>, 2024, accessed: 2024-06-18.
- [5] D. Guo, Q. Zhu, D. Yang, Z. Xie, K. Dong, W. Zhang, G. Chen, X. Bi, Y. Wu, Y. Li *et al.*, “Deepseek-coder: When the large language model meets programming—the rise of code intelligence,” *arXiv preprint arXiv:2401.14196*, 2024.
- [6] Cursor, “Cursor: Ai-powered coding assistant,” <https://www.cursor.com/>, 2024, accessed: 2024-06-18.
- [7] D. Nam, A. Macvean, V. Hellendoorn, B. Vasilescu, and B. Myers, “Using an llm to help with code understanding,” in *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*, 2024, pp. 1–13.
- [8] J. Bai, S. Bai, Y. Chu, Z. Cui, K. Dang, X. Deng, Y. Fan, W. Ge, Y. Han, F. Huang *et al.*, “Qwen technical report,” *arXiv preprint arXiv:2309.16609*, 2023.
- [9] F. Zhang, B. Chen, Y. Zhang, J. Keung, J. Liu, D. Zan, Y. Mao, J.-G. Lou, and W. Chen, “Repocoder: Repository-level code completion through iterative retrieval and generation,” in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 2023, pp. 2471–2484.
- [10] Z. Li, C. Wang, Z. Liu, H. Wang, D. Chen, S. Wang, and C. Gao, “Cctest: Testing and repairing code completion systems,” in *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 2023, pp. 1238–1250.
- [11] P. Nie, R. Banerjee, J. J. Li, R. J. Mooney, and M. Gligoric, “Learning deep semantics for test completion,” in *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 2023, pp. 2111–2123.
- [12] T. Ahmed and P. Devanbu, “Few-shot training llms for project-specific code-summarization,” in *Proceedings of the 37th IEEE/ACM International Conference on Automated Software Engineering*, 2022, pp. 1–5.
- [13] D. Sobania, M. Briesch, C. Hanna, and J. Petke, “An analysis of the automatic bug fixing performance of chatgpt,” in *2023 IEEE/ACM International Workshop on Automated Program Repair (APR)*. IEEE, 2023, pp. 23–30.
- [14] Z. Ni, B. Li, X. Sun, T. Chen, B. Tang, and X. Shi, “Analyzing bug fix for automatic bug cause classification,” *Journal of Systems and Software*, vol. 163, p. 110538, 2020.
- [15] D. Drain, C. Wu, A. Svyatkovskiy, and N. Sundaresan, “Generating bug-fixes using pretrained transformers,” in *Proceedings of the 5th ACM SIGPLAN International Symposium on Machine Programming*, 2021, pp. 1–8.
- [16] GitHub, “Github copilot,” <https://github.com/features/copilot>, 2023, accessed: 2024-06-19.
- [17] JetBrains, “Jetbrains: Essential tools for software developers and teams,” <https://www.jetbrains.com/>, 2024, accessed: 2024-06-19.
- [18] H. Yu, B. Shen, D. Ran, J. Zhang, Q. Zhang, Y. Ma, G. Liang, Y. Li, Q. Wang, and T. Xie, “Codereval: A benchmark of pragmatic code generation with generative pre-trained models,” in *Proceedings of the 46th IEEE/ACM International Conference on Software Engineering*, 2024, pp. 1–12.
- [19] K. Zhang, J. Li, G. Li, X. Shi, and Z. Jin, “Codeagent: Enhancing code generation with tool-integrated agent systems for real-world repo-level coding challenges,” *arXiv preprint arXiv:2401.07339*, 2024.
- [20] M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. d. O. Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman *et al.*, “Evaluating large language models trained on code,” *arXiv preprint arXiv:2107.03374*, 2021.
- [21] Y. Ding, Z. Wang, W. U. Ahmad, H. Ding, M. Tan, N. Jain, M. K. Ramanathan, R. Nallapati, P. Bhatia, D. Roth, and B. Xiang, “Crosscodeeval: A diverse and multilingual benchmark for cross-file code completion,” in *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. [Online]. Available: <https://openreview.net/forum?id=wgDcbBMSfh>
- [22] C. E. Jimenez, J. Yang, A. Wettig, S. Yao, K. Pei, O. Press, and K. R. Narasimhan, “Swe-bench: Can language models resolve real-world github issues?” in *The Twelfth International Conference on Learning Representations*, 2023.
- [23] S. P. Sahu, M. Mandal, S. Bharadwaj, A. Kanade, P. Maniatis, and S. Shevade, “Codequeries: A dataset of semantic queries over code,” in *Proceedings of the 17th Innovations in Software Engineering Conference*, 2024, pp. 1–11.
- [24] C. Liu and X. Wan, “Codeqa: A question answering dataset for source code comprehension,” in *Findings of the Association for Computational Linguistics: EMNLP 2021*, 2021, pp. 2618–2632.
- [25] C. Lee, Y. Seonwoo, and A. Oh, “Cs1qa: A dataset for assisting code-based question answering in an introductory programming course,” in *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2022, pp. 2026–2040.
- [26] A. Bansal, Z. Eberhart, L. Wu, and C. McMillan, “A neural question answering system for basic questions about subroutines,” in *2021 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER)*. IEEE, 2021, pp. 60–71.
- [27] L. Zheng, W.-L. Chiang, Y. Sheng, S. Zhuang, Z. Wu, Y. Zhuang, Z. Lin, Z. Li, D. Li, E. Xing, H. Zhang, J. E. Gonzalez, and I. Stoica, “Judging llm-as-a-judge with mt-bench and chatbot arena,” in *Advances in Neural Information Processing Systems*, A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, Eds., vol. 36. Curran Associates, Inc., 2023, pp. 46 595–46 623.
- [28] S. Es, J. James, L. Espinosa-Anke, and S. Schockaert, “Ragas: Automated evaluation of retrieval augmented generation,” *arXiv preprint arXiv:2309.15217*, 2023.
- [29] S. Robertson, H. Zaragoza, and M. Taylor, “Simple bm25 extension to multiple weighted fields,” in *Proceedings of the thirteenth ACM international conference on Information and knowledge management*, 2004, pp. 42–49.
- [30] T. Ben-Nun, A. S. Jakobovits, and T. Hoeffler, “Neural code comprehension: A learnable representation of code semantics,” *Advances in neural information processing systems*, vol. 31, 2018.
- [31] Z. Porkoláb, T. Brunner, D. Krupp, and M. Csordás, “Codecompass: an open software comprehension framework for industrial usage,” in *Proceedings of the 26th Conference on Program Comprehension*, 2018, pp. 361–369.
- [32] F. Bertolotti and W. Cazzola, “Fold2vec: Towards a statement-based representation of code for code comprehension,” *ACM Transactions on Software Engineering and Methodology*, vol. 32, no. 1, pp. 1–31, 2023.
- [33] Q. Luo, Y. Ye, S. Liang, Z. Zhang, Y. Qin, Y. Lu, Y. Wu, X. Cong, Y. Lin, Y. Zhang *et al.*, “Repoagent: An llm-powered open-source framework for repository-level code documentation generation,” *arXiv preprint arXiv:2402.16667*, 2024.
- [34] P. Jayarao and A. Srivastava, “Intent detection for code-mix utterances in task oriented dialogue systems,” in *2018 International Conference on Electrical, Electronics, Communication, Computer, and Optimization Techniques (ICEECCOT)*, 2018, pp. 583–587.
- [35] K. Pearce, S. Alghowinem, and C. Breazeal, “Build-a-bot: teaching conversational ai using a transformer-based intent recognition and question answering architecture,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 13, 2023, pp. 16 025–16 032.
- [36] Q. Dong, L. Li, D. Dai, C. Zheng, Z. Wu, B. Chang, X. Sun, J. Xu, and Z. Sui, “A survey on in-context learning,” *arXiv preprint arXiv:2301.00234*, 2022.
- [37] Y. Ren, Y. Cao, P. Guo, F. Fang, W. Ma, and Z. Lin, “Retrieve-and-sample: Document-level event argument extraction via hybrid retrieval augmentation,” in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2023, pp. 293–306.

- 
- [38] X. Li, J. Jin, Y. Zhou, Y. Zhang, P. Zhang, Y. Zhu, and Z. Dou, "From matching to generation: A survey on generative information retrieval," *arXiv preprint arXiv:2404.14851*, 2024.
- [39] M. Poli, S. Massaroli, E. Nguyen, D. Y. Fu, T. Dao, S. Baccus, Y. Bengio, S. Ermon, and C. Ré, "Hyena hierarchy: Towards larger convolutional language models," in *International Conference on Machine Learning*. PMLR, 2023, pp. 28 043–28 078.
- [40] Z. Dai, Z. Yang, Y. Yang, J. Carbonell, Q. V. Le, and R. Salakhutdinov, "Transformer-xl: Attentive language models beyond a fixed-length context," *arXiv preprint arXiv:1901.02860*, 2019.
- [41] P. Zhang, Z. Liu, S. Xiao, N. Shao, Q. Ye, and Z. Dou, "Soaring from 4k to 400k: Extending llm's context with activation beacon," *arXiv preprint arXiv:2401.03462*, 2024.
- [42] Hugging Face, "The Stack: A Large-scale Context-Aware Code Dataset," <https://huggingface.co/datasets/bigcode/the-stack>, Jan. 2023.
- [43] OpenAI, "Openai tokenizer," <https://platform.openai.com/tokenizer>, 2024, accessed: 2024-06-18.
- [44] X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, S. Narang, A. Chowdhery, and D. Zhou, "Self-consistency improves chain of thought reasoning in language models," *arXiv preprint arXiv:2203.11171*, 2022.
- [45] T. Fredriksson, D. Issa Mattos, J. Bosch, and H. Olsson, *Data Labeling: An Empirical Investigation into Industrial Challenges and Mitigation Strategies*, 11 2020, pp. 202–216.
- [46] H. Chase, "LangChain," Oct. 2022. [Online]. Available: <https://github.com/langchain-ai/langchain>
- [47] D. Brown, "Rank-BM25: A Collection of BM25 Algorithms in Python," 2020. [Online]. Available: <https://doi.org/10.5281/zenodo.4520057>
- [48] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou *et al.*, "Chain-of-thought prompting elicits reasoning in large language models," *Advances in neural information processing systems*, vol. 35, pp. 24 824–24 837, 2022.
- [49] T. Li, W.-L. Chiang, E. Frick, L. Dunlap, T. Wu, B. Zhu, J. E. Gonzalez, and I. Stoica, "From crowdsourced data to high-quality benchmarks: Arena-hard and benchbuilder pipeline," *arXiv preprint arXiv:2406.11939*, 2024.
- [50] W.-L. Chiang, L. Zheng, Y. Sheng, A. N. Angelopoulos, T. Li, D. Li, H. Zhang, B. Zhu, M. Jordan, J. E. Gonzalez, and I. Stoica, "Chatbot arena: An open platform for evaluating llms by human preference," 2024.
- [51] T. Liu, C. Xu, and J. McAuley, "Repobench: Benchmarking repository-level code auto-completion systems," *arXiv:2306.03091*, Jun 2023.
- [52] OpenAI, "Six strategies for getting better results," <https://platform.openai.com/docs/guides/prompt-engineering/six-strategies-for-getting-better-results>, 2024.
- [53] BigCode, "The stack dataset," 2024, collection of source code in over 300 programming languages. [Online]. Available: <https://huggingface.co/datasets/bigcode/the-stack-dedup/blob/main/licenses.json>